

# Deep Learning: Overview and basics

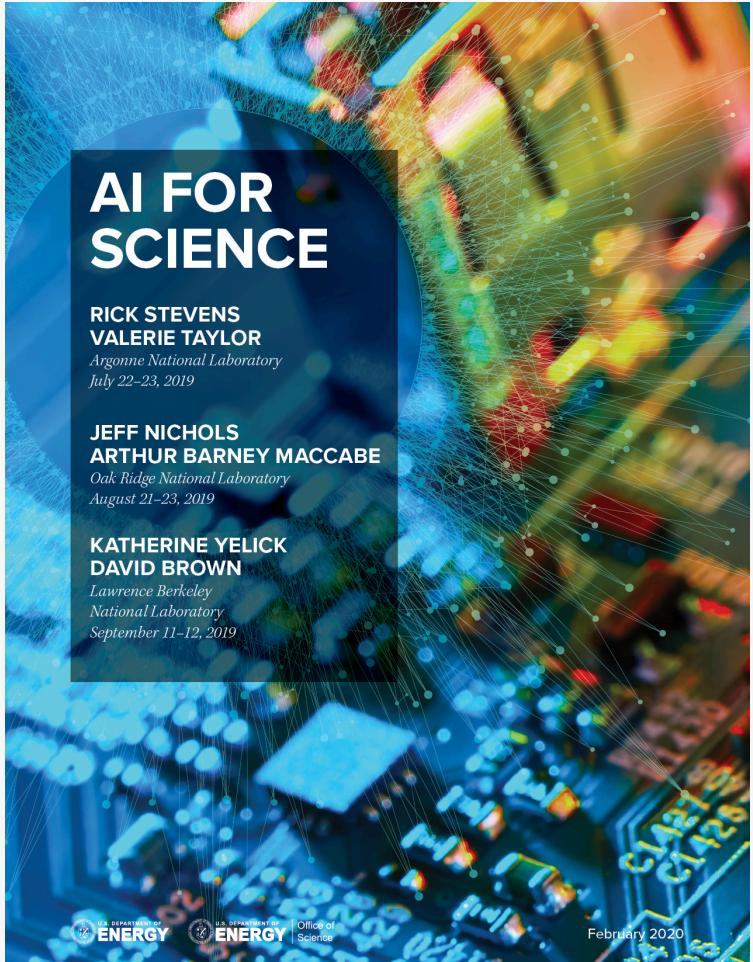
ATPESC 2020  
Friday 2020-03-19

Kyle Gerard Felker  
Postdoctoral Appointee, Aurora Early Science Program  
Argonne Leadership Computing Facility  
[felker@anl.gov](mailto:felker@anl.gov)

# Outline

- **High-level introduction to deep learning (25 min)**
- **Hands-on exercises (20 min)**

# AI for Science



**NERSC: Perlmutter**  
AMD Epyc + NVIDIA A100s  
(2020)



**ALCF: Aurora**  
Intel Xeon + Xe GPUs  
(2021)



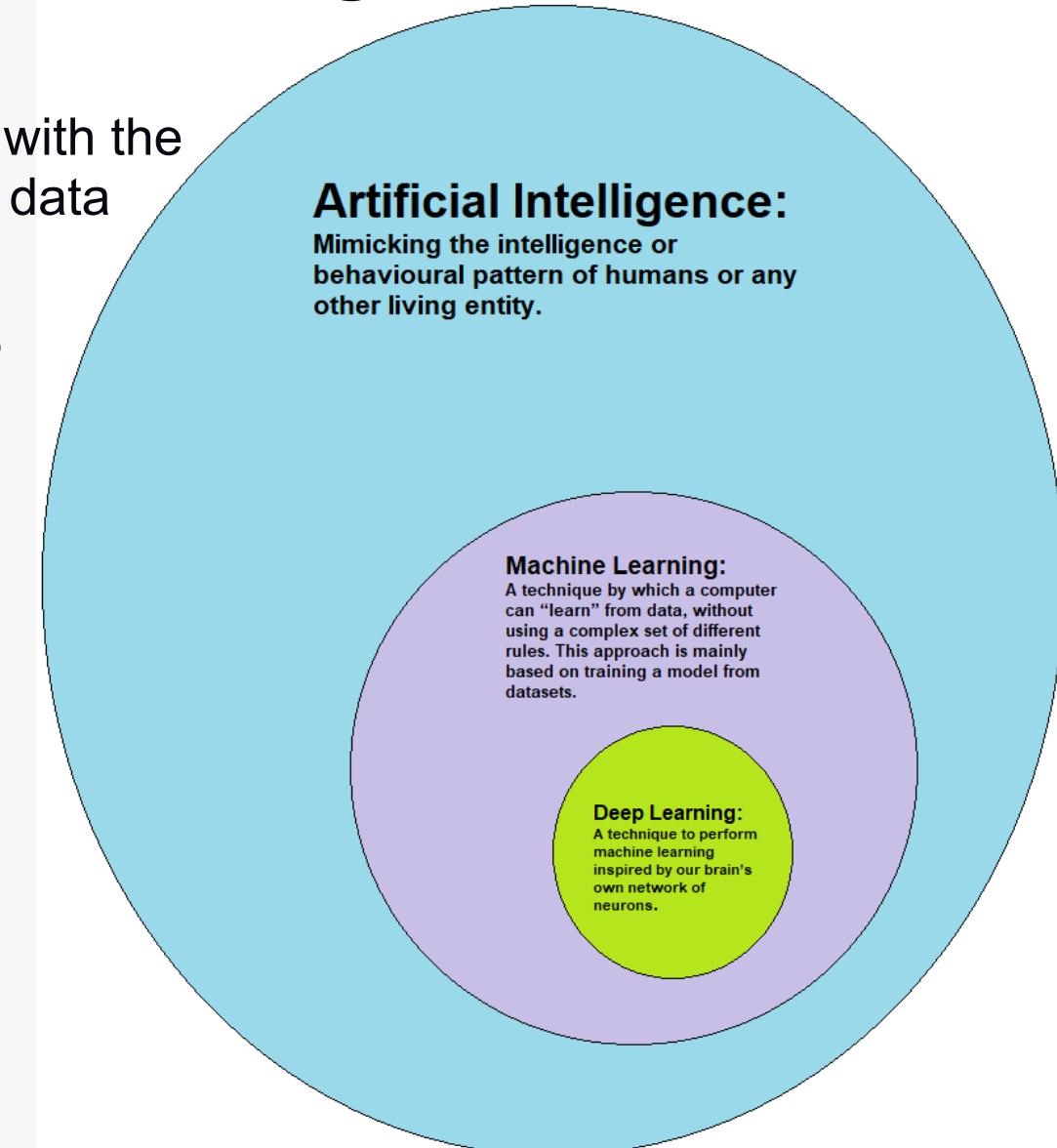
**NERSC: Perlmutter**  
AMD Epyc + GPUs  
(2021)

# AI, machine learning, and deep learning

“Machine Learning is a subfield of artificial intelligence with the goal of developing algorithms capable of learning from data automatically.”

Mehta, et al (2019)

“A high-bias, low-variance introduction to Machine Learning for physicists”



# AI, machine learning, and deep learning

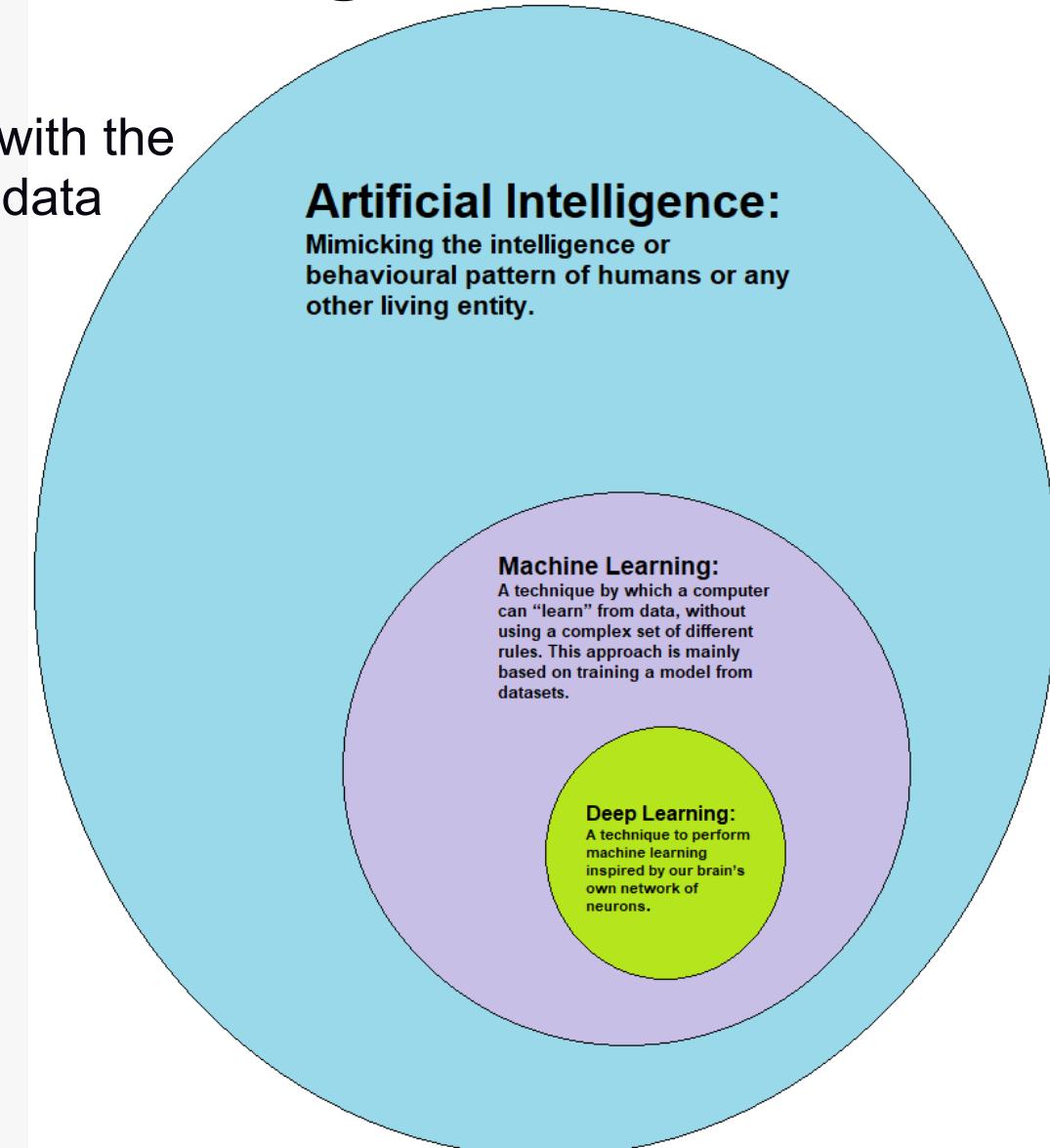
“Machine Learning is a subfield of artificial intelligence with the goal of developing algorithms capable of learning from data automatically.”

Mehta, et al (2019)

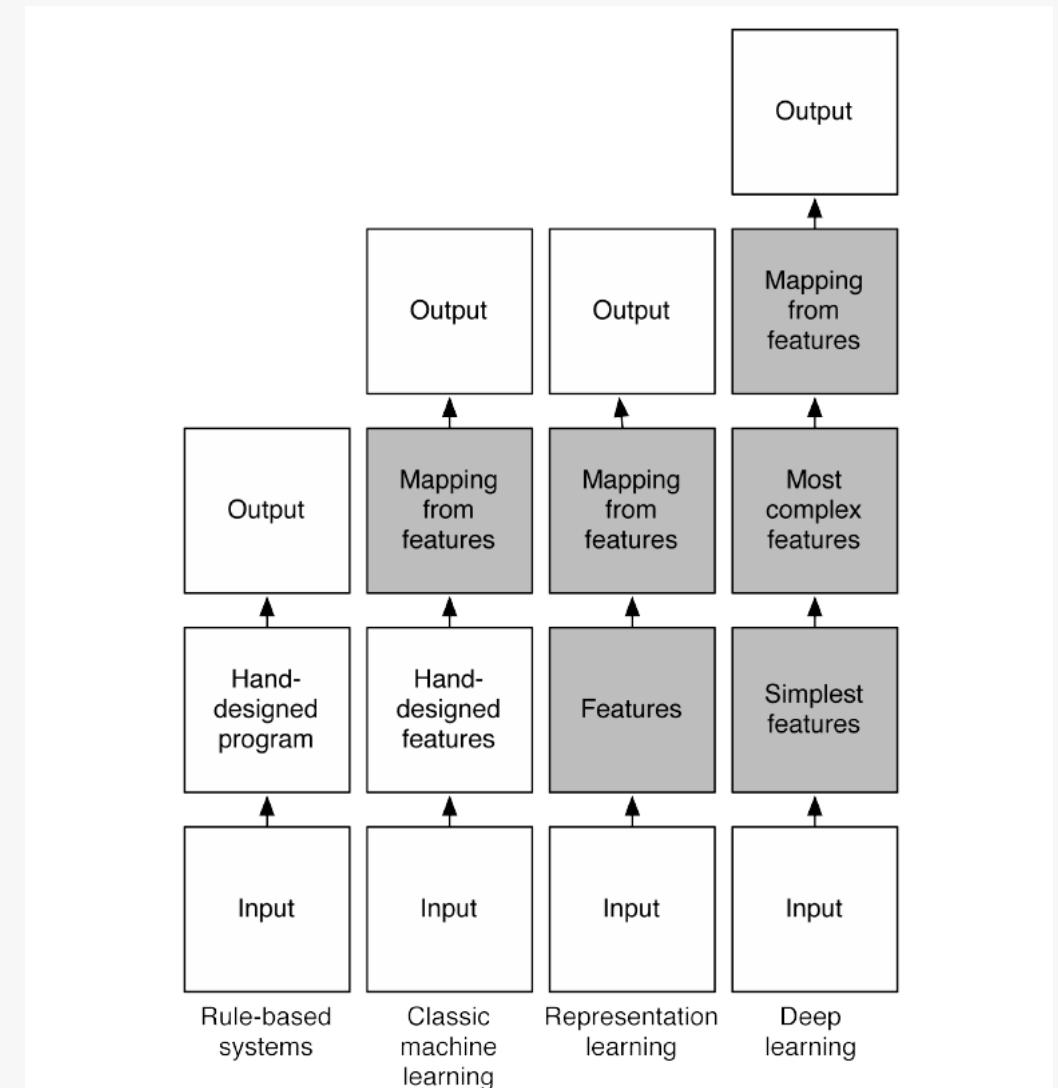
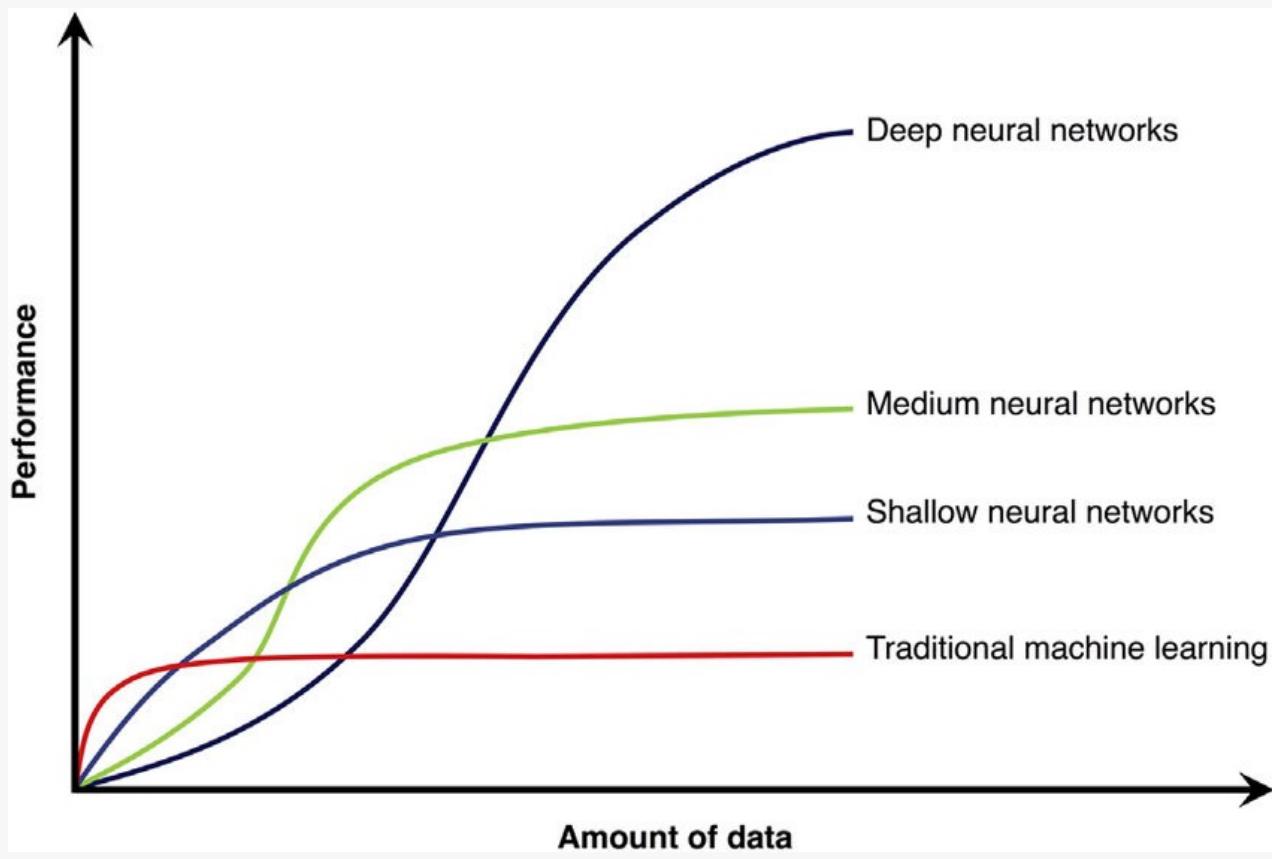
“A high-bias, low-variance introduction to Machine Learning for physicists”

## Machine learning algorithms:

- can improve their accuracy given more data
- are capable of “teaching themselves”
- gain knowledge that is not programmed with an explicit set of rules



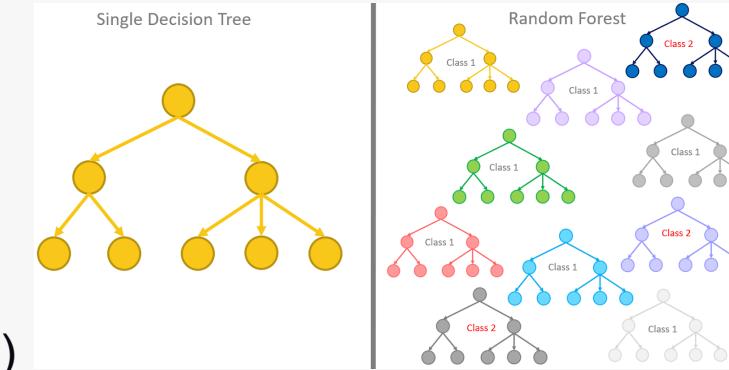
# Motivations for deep learning



# Categories of learning problems or paradigms

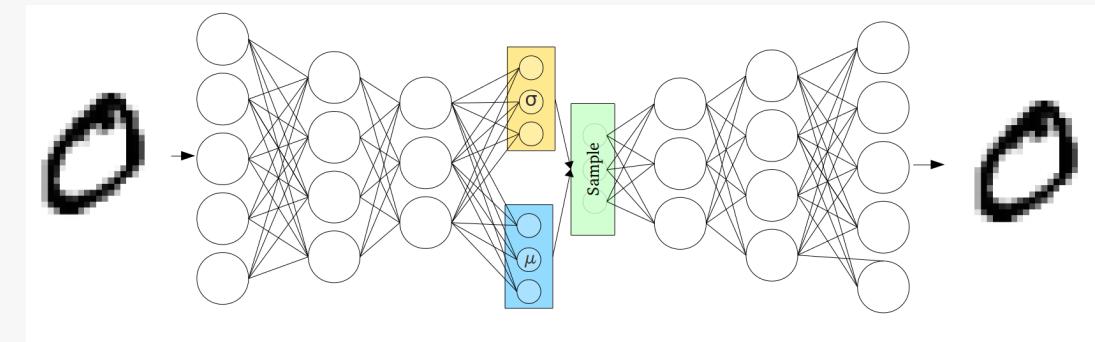
- Supervised learning (**this talk**)
  - Regression: output variable is continuous
  - Classification: output variable is discrete (categorical)
- Unsupervised learning (**Corey and Romit's talk**)
  - Clustering
  - Association
- Semi-supervised learning
  - *Mostly* unlabeled data
- Reinforcement learning

## Examples of algorithms



**Supervised:** decision tree and random forests

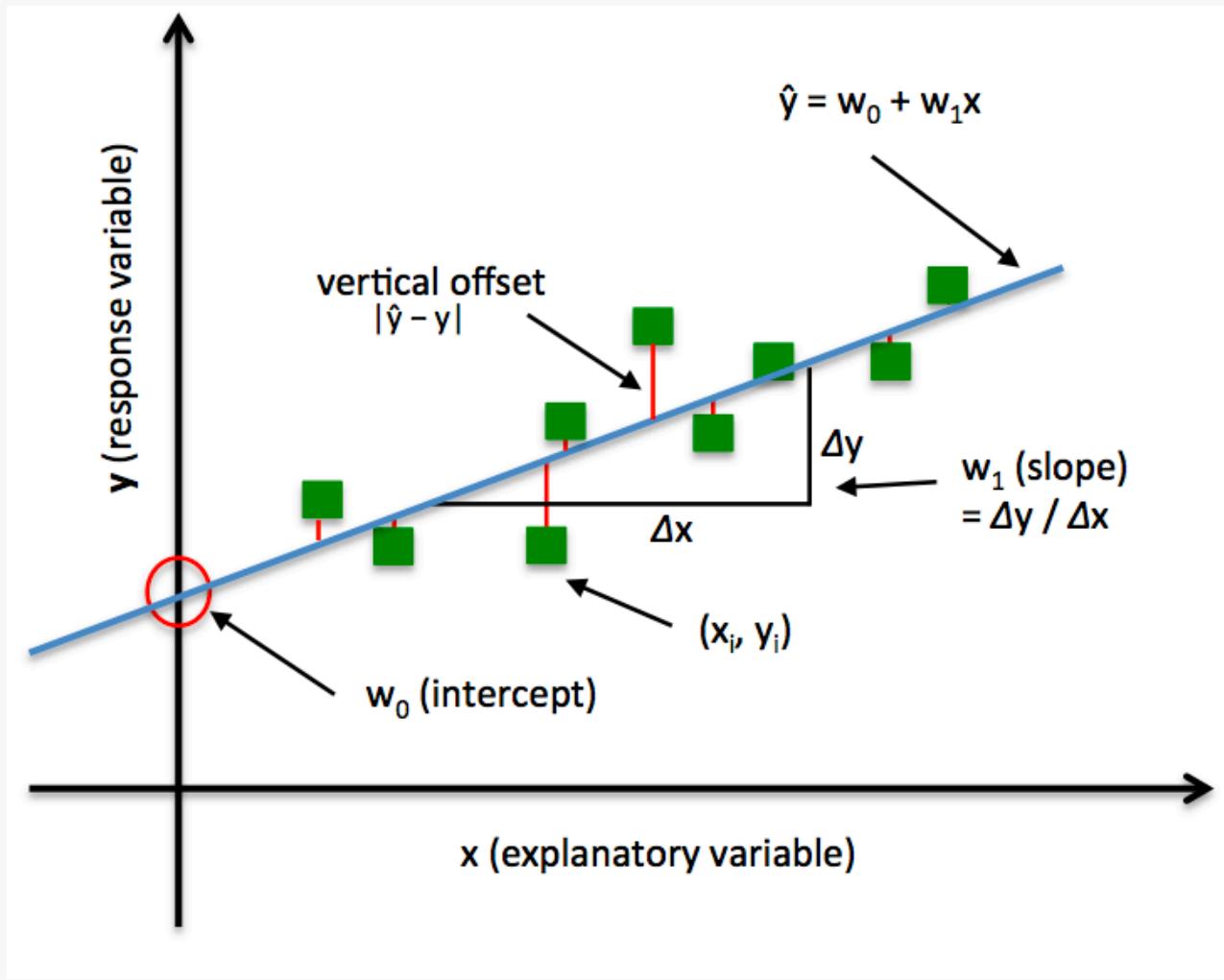
<https://towardsdatascience.com/from-a-single-decision-tree-to-a-random-forest-b9523be65147>



**Unsupervised:** Variational Autoencoder (VAE)

<https://towardsdatascience.com/intuitively-understanding-variational-autoencoders-1bfe67eb5daf>

# Linear regression: ordinary least squares (OLS)



<https://sebastianraschka.com/faq/docs/closed-form-vs-gd.html>

Minimizes the Mean Squared Error (MSE), or **quadratic loss**:

$$J = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

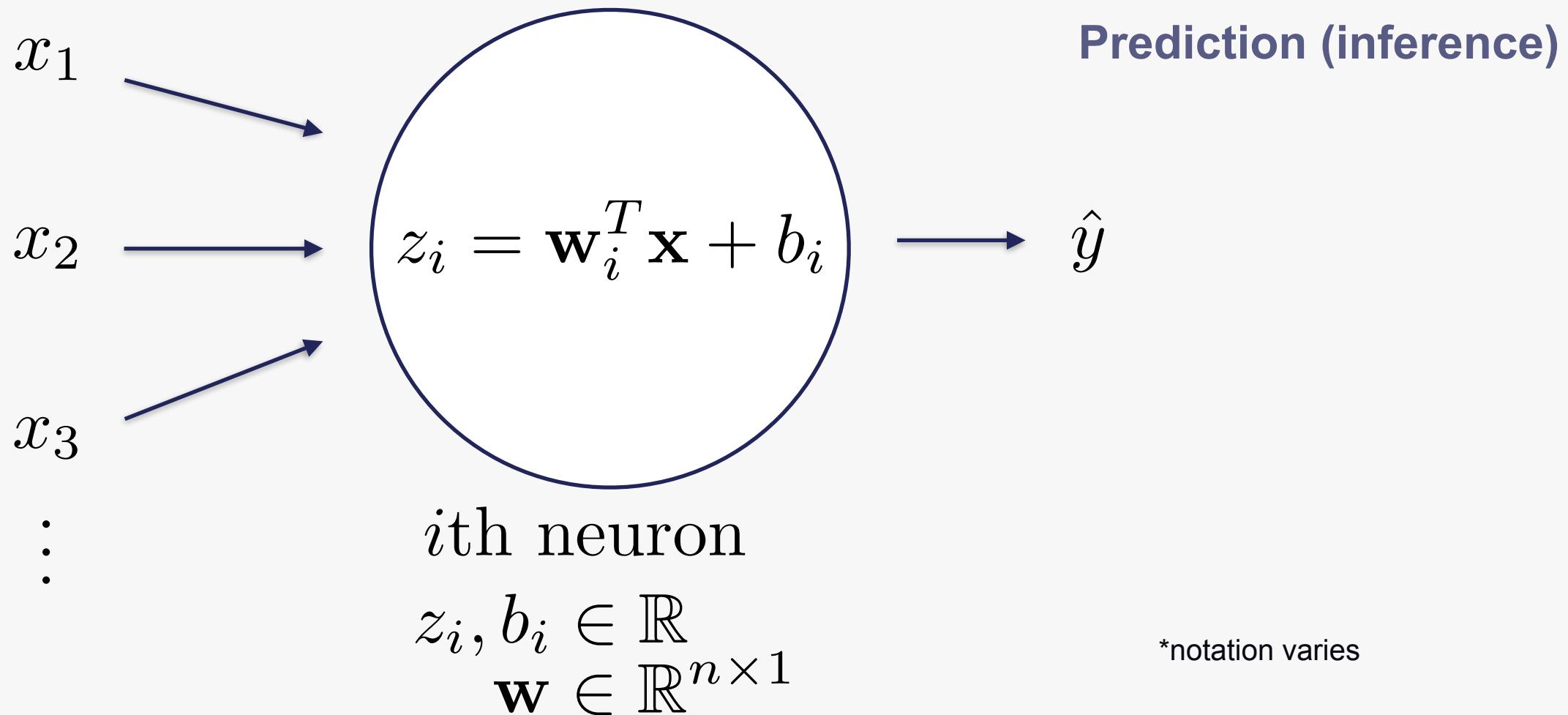
Closed-form solution exists;  
one-step procedure: **normal equation**

$$\mathbf{w}_i = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Alternative solutions may be required or preferred in some cases

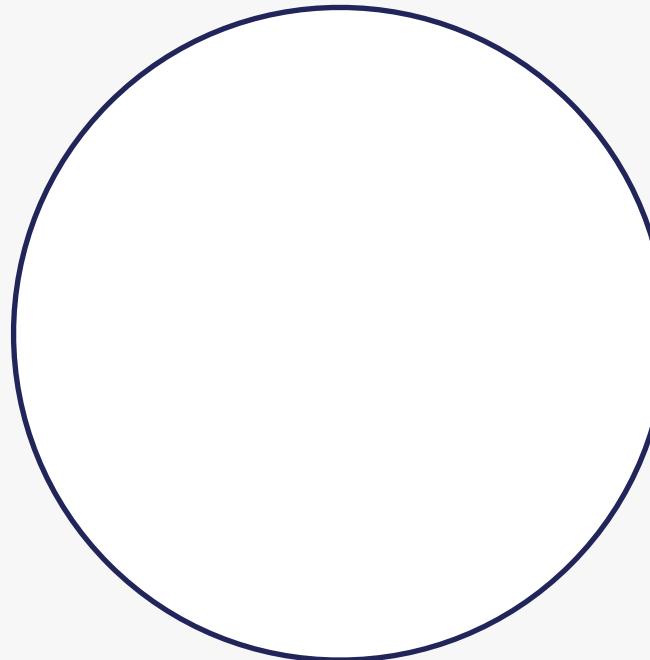
# Neuron (linear regression)

Input features  $\mathbf{x} \in \mathbb{R}^{n \times 1}$



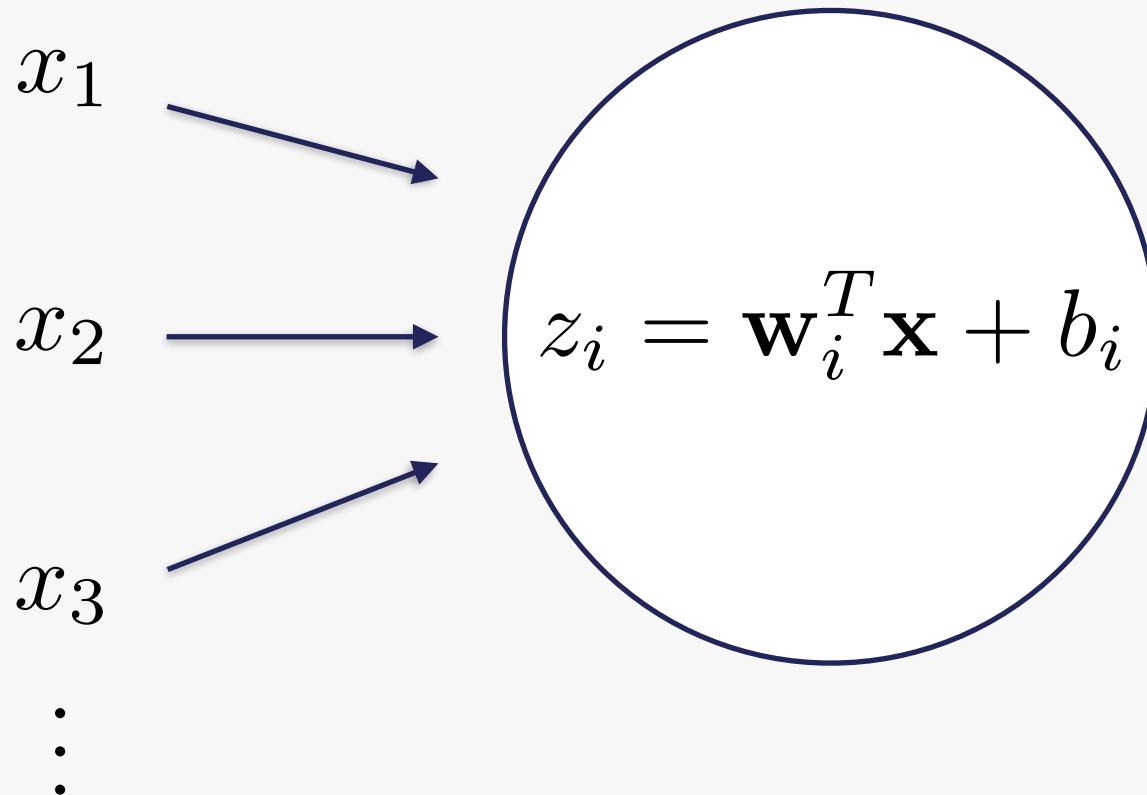
# Forward pass: compute the loss

Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



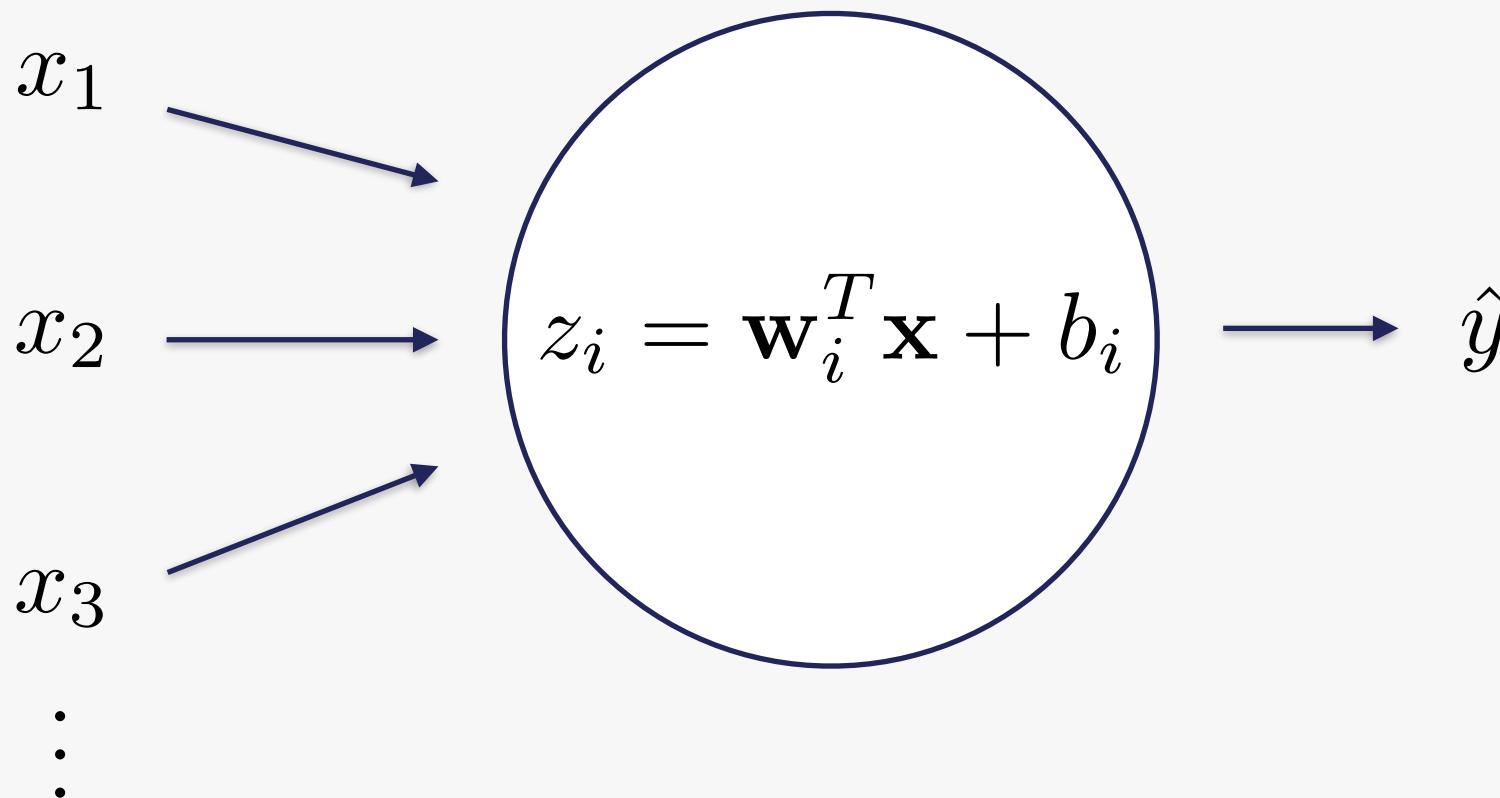
# Forward pass: compute the loss

Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



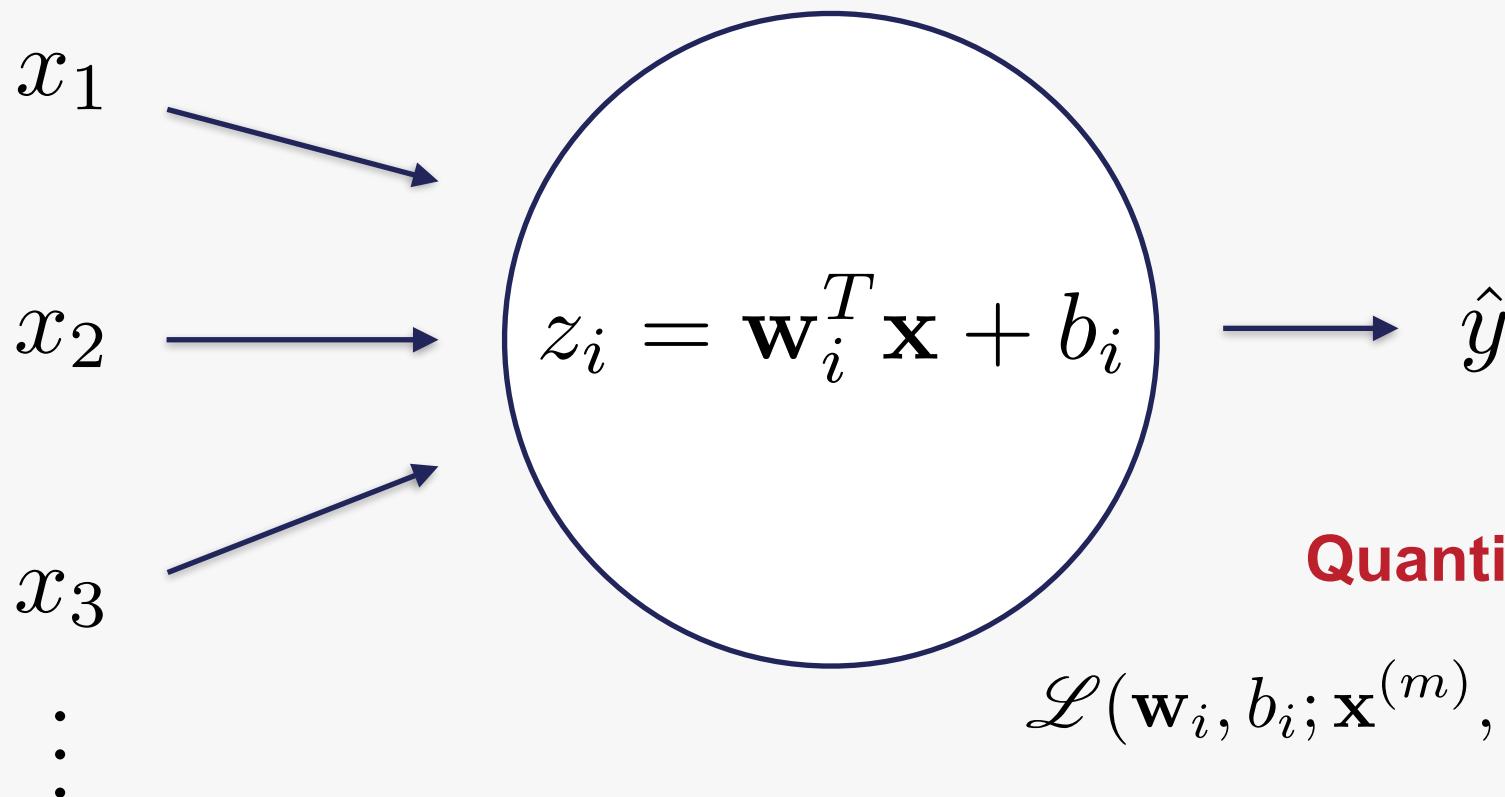
# Forward pass: compute the loss

Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



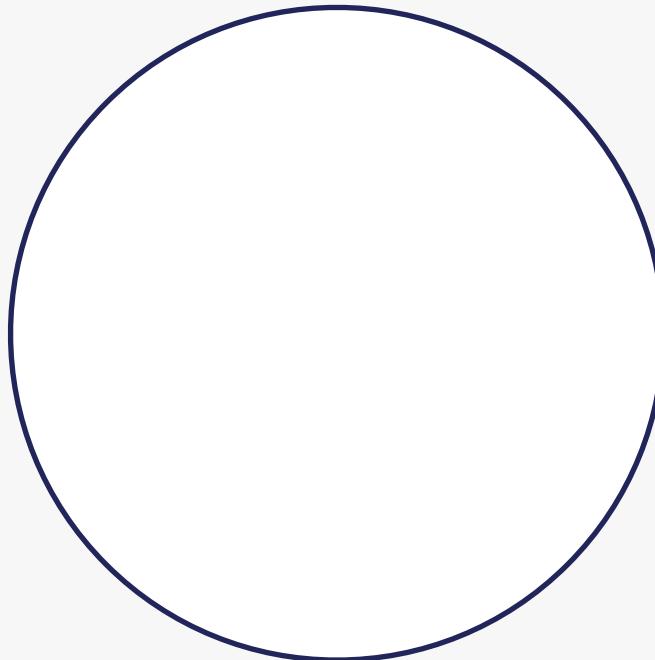
# Forward pass: compute the loss

Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



# Backward pass: compute the derivatives

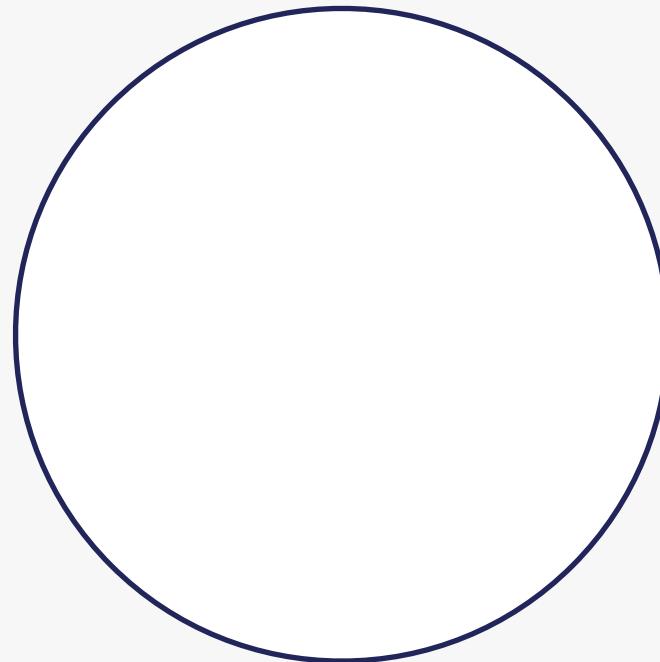
Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



$$\mathcal{L}(\mathbf{w}_i, b_i; \mathbf{x}^{(m)}, y^{(m)}) = (\hat{y} - y^{(m)})^2$$

# Backward pass: compute the derivatives

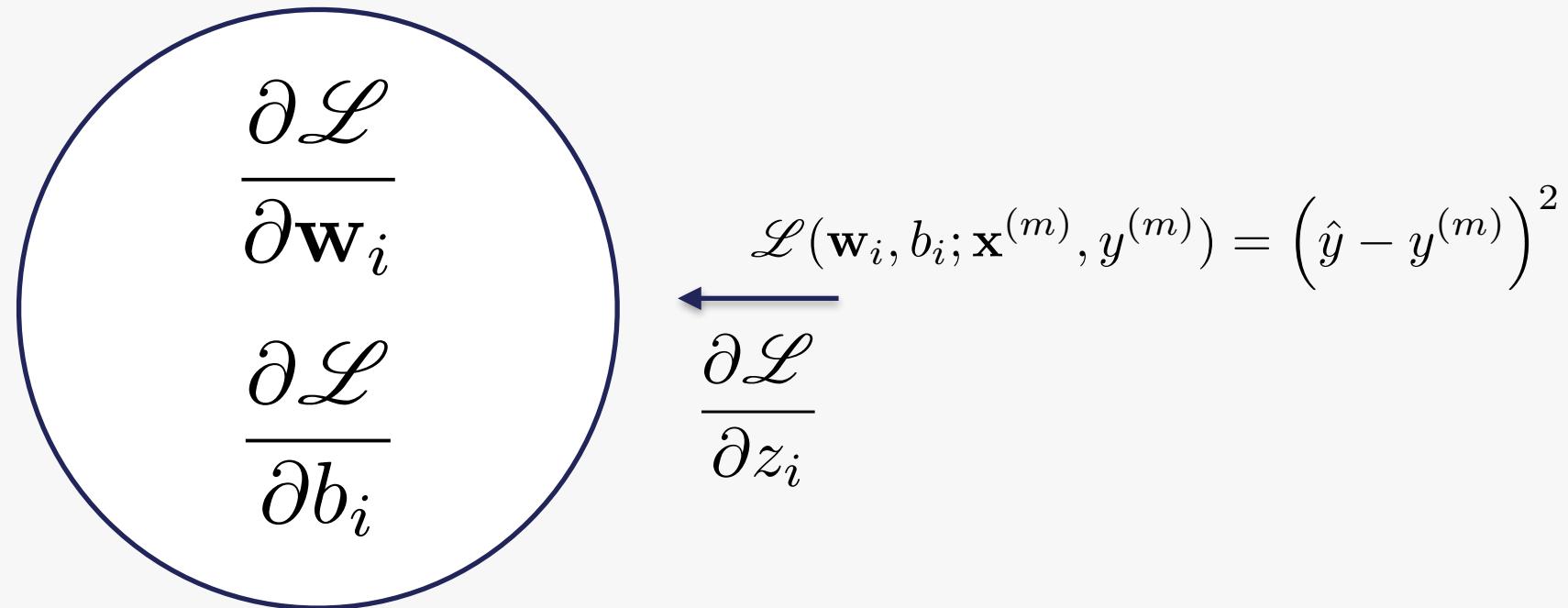
Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



$$\frac{\partial \mathcal{L}}{\partial z_i} \leftarrow \mathcal{L}(\mathbf{w}_i, b_i; \mathbf{x}^{(m)}, y^{(m)}) = (\hat{y} - y^{(m)})^2$$

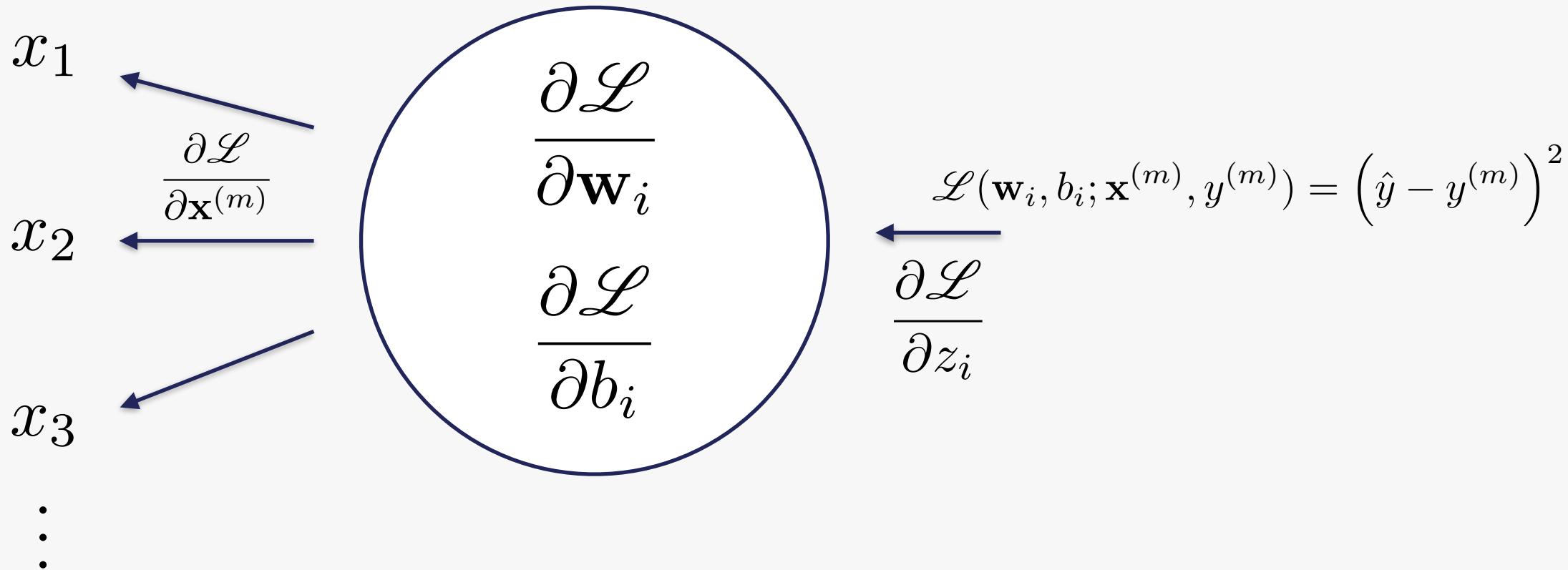
# Backward pass: compute the derivatives

Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



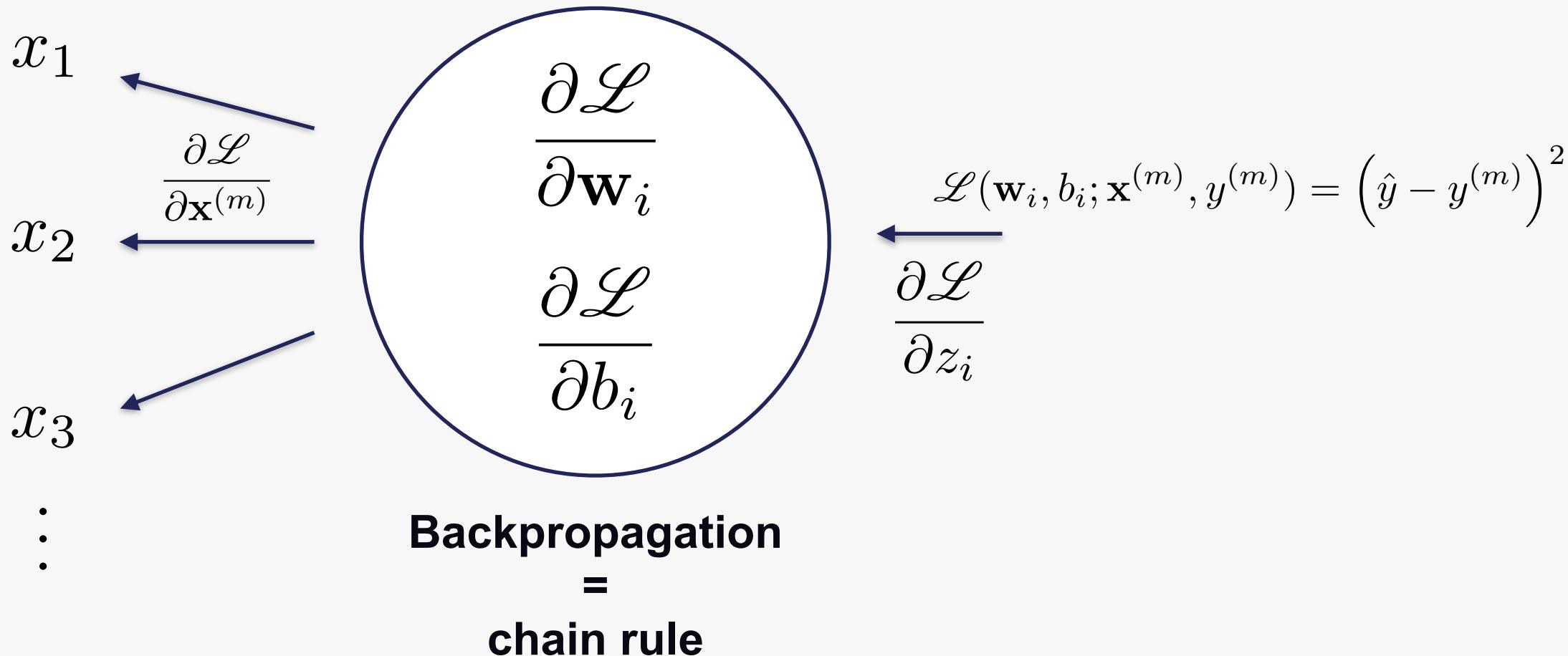
# Backward pass: compute the derivatives

Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



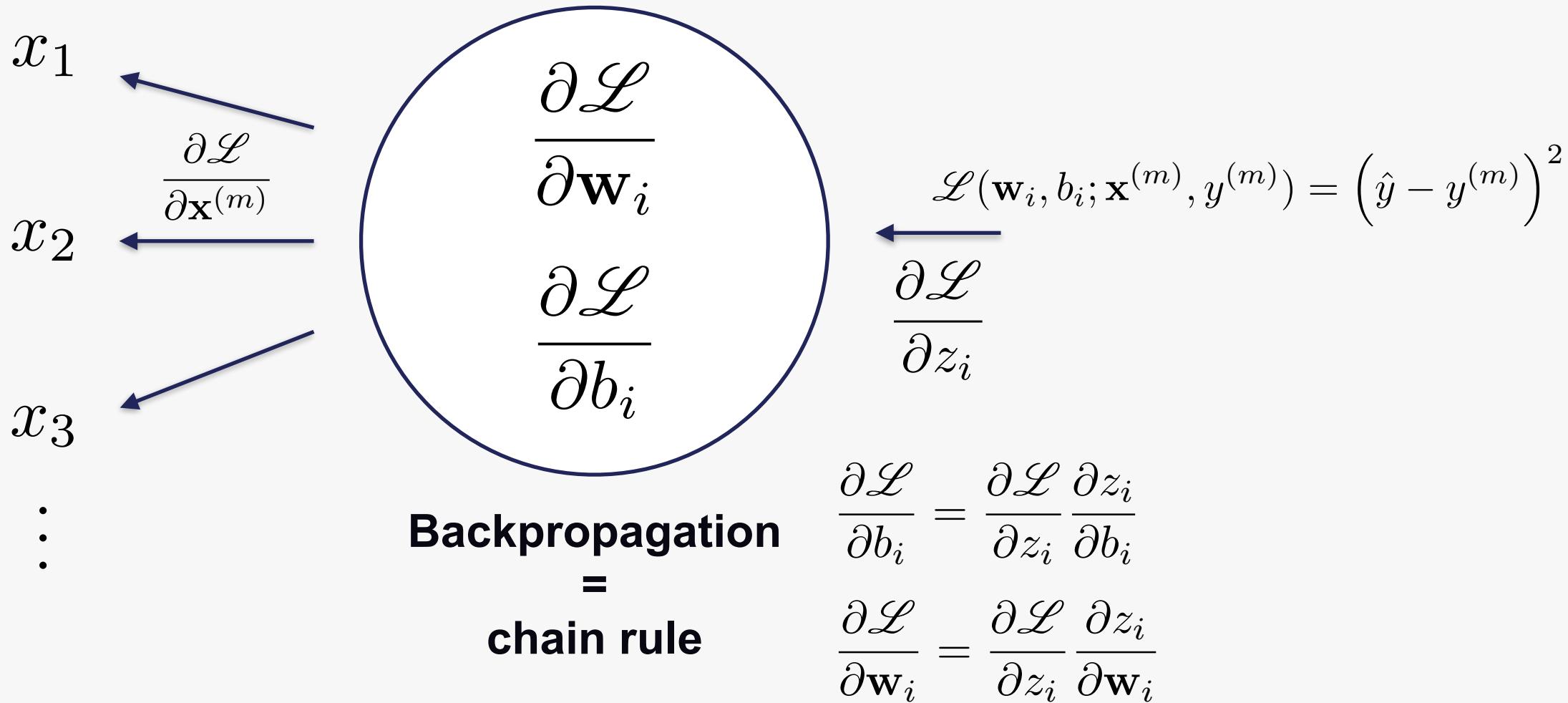
# Backward pass: compute the derivatives

Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



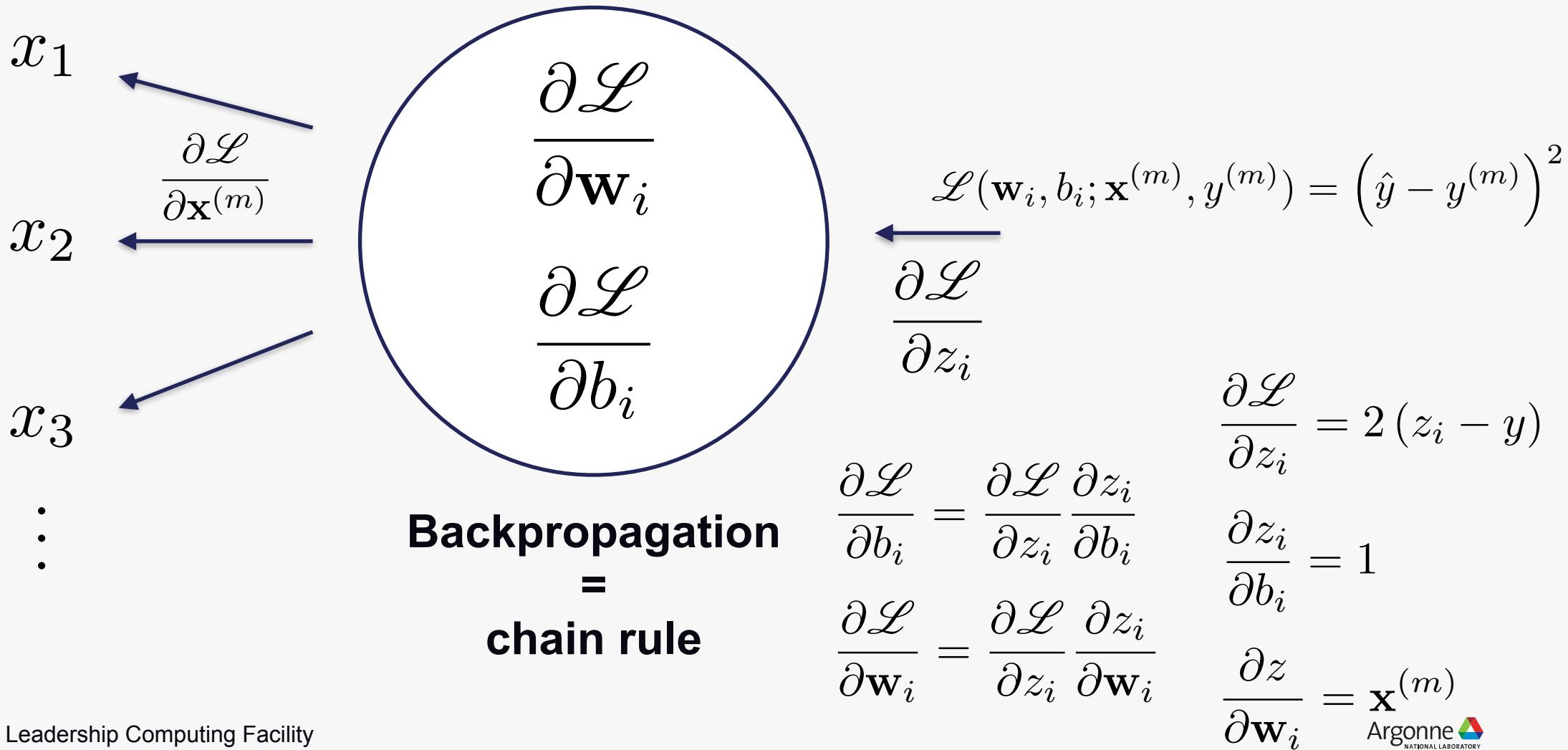
# Backward pass: compute the derivatives

Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



# Backward pass: compute the derivatives

Given single labeled training example,  $(\mathbf{x}^{(m)}, y^{(m)})$



# Training neuron weights with gradient descent

Let  $\theta \equiv \{\mathbf{w}_i, b_i\}$  trainable parameters

**Cost:** average loss over training examples

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\theta; \mathbf{x}^{(m)}, y^{(m)})$$

- Initialize  $\theta$  randomly
- For N epochs
  - For each batch of training examples  $\{(x_0, y_0), \dots, (x_b, y_b)\}$ :
    - \* compute loss gradient:  $\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{B} \sum_{i=1}^N \frac{\partial J^i(\theta)}{\partial \theta}$
    - \* update  $\theta$  with update rule:
$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$

# Training neuron weights with gradient descent

Let  $\theta \equiv \{\mathbf{w}_i, b_i\}$  trainable parameters

**Cost:** average loss over training examples

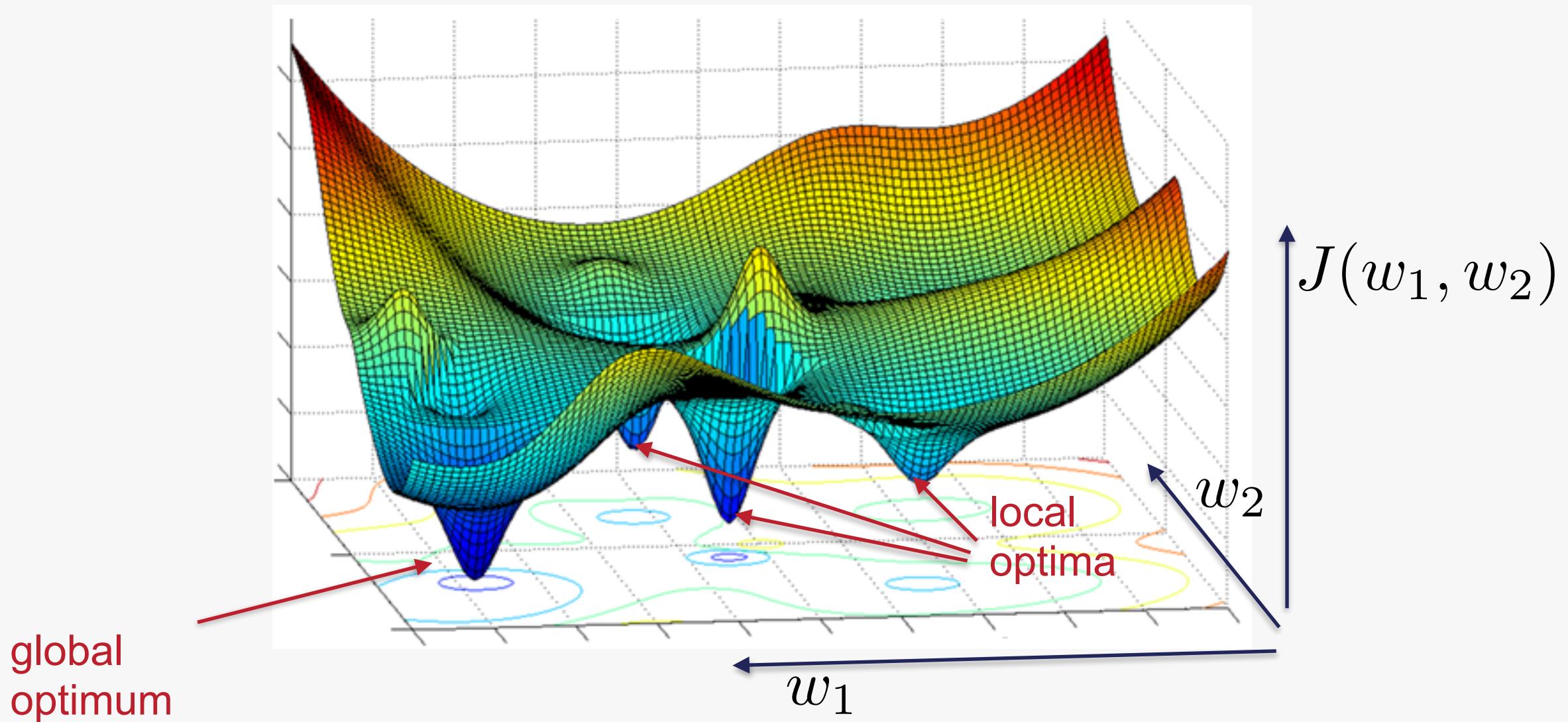
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\theta; \mathbf{x}^{(m)}, y^{(m)})$$

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$$

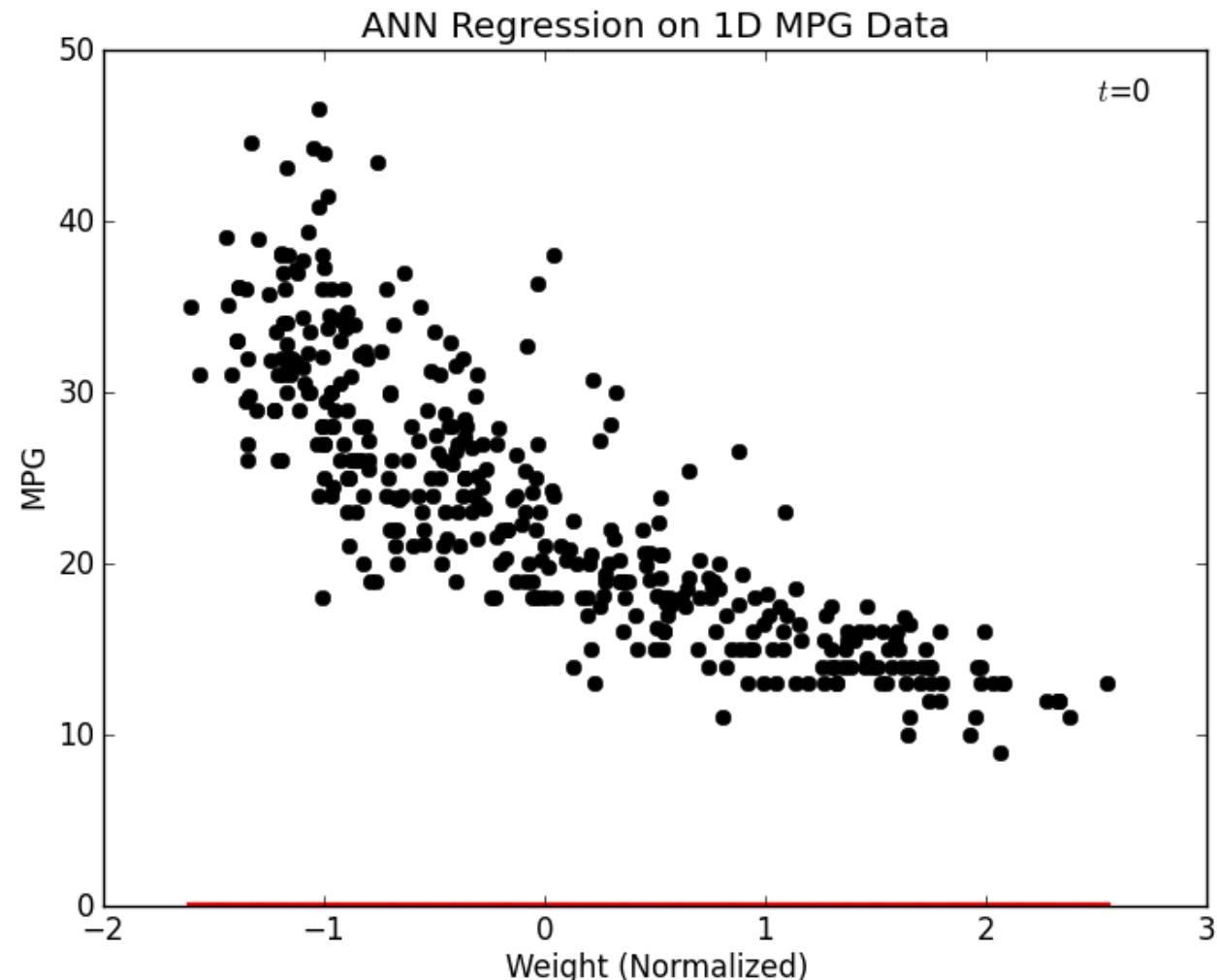
$\eta \equiv$  learning rate

- Initialize  $\theta$  randomly
- For N epochs
  - For each batch of training examples  $\{(x_0, y_0), \dots, (x_b, y_b)\}$ :
    - \* compute loss gradient:  $\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{B} \sum_{i=1}^N \frac{\partial J^i(\theta)}{\partial \theta}$
    - \* update  $\theta$  with update rule:  
$$\theta := \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$$

# Global vs. local minima



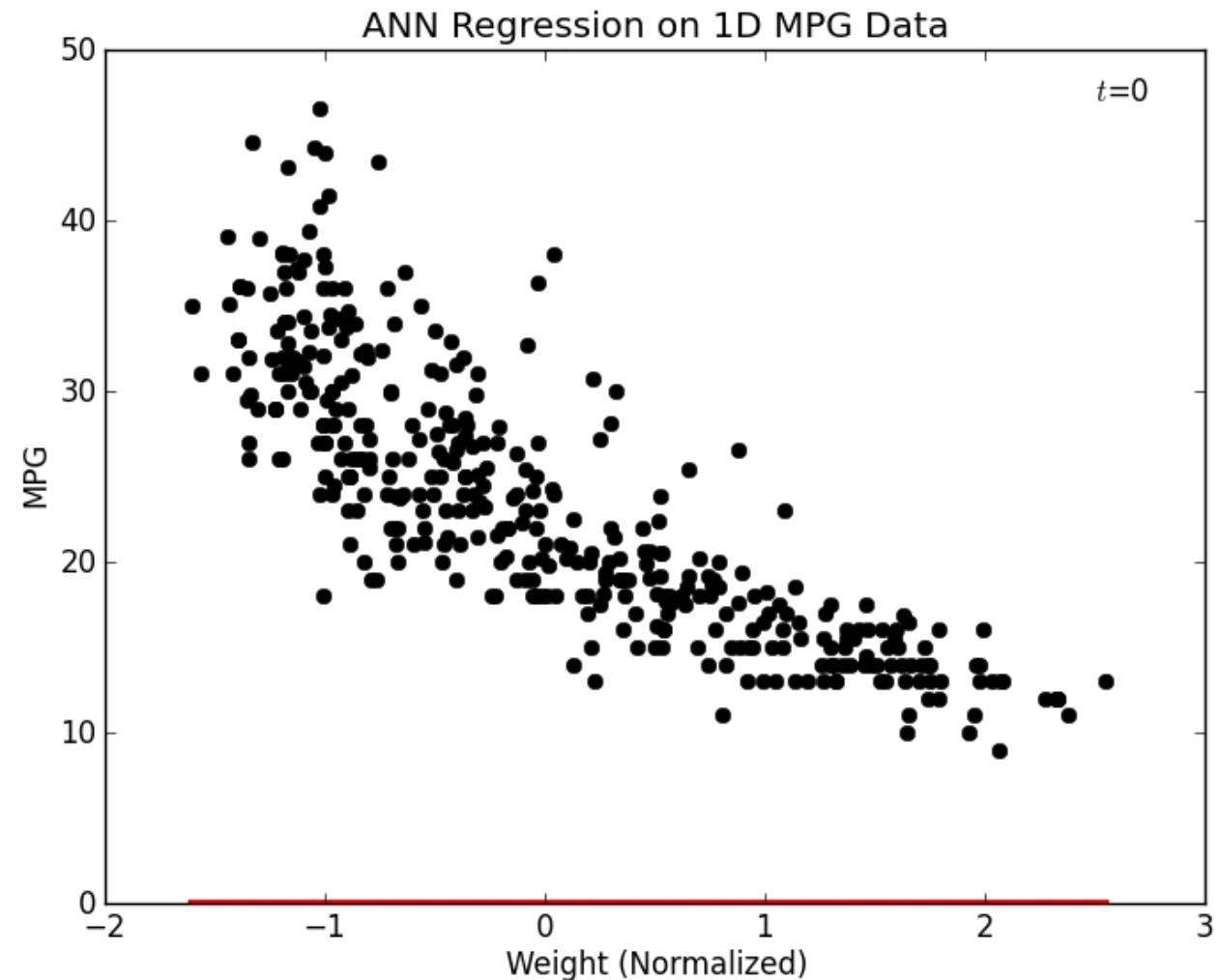
# OLS solved with gradient descent



$t = \text{epoch}$   
= neuron sees  
**all samples in the**  
training set

<http://www.briandolhansky.com/blog/artificial-neural-networks-linear-regression-part-1>

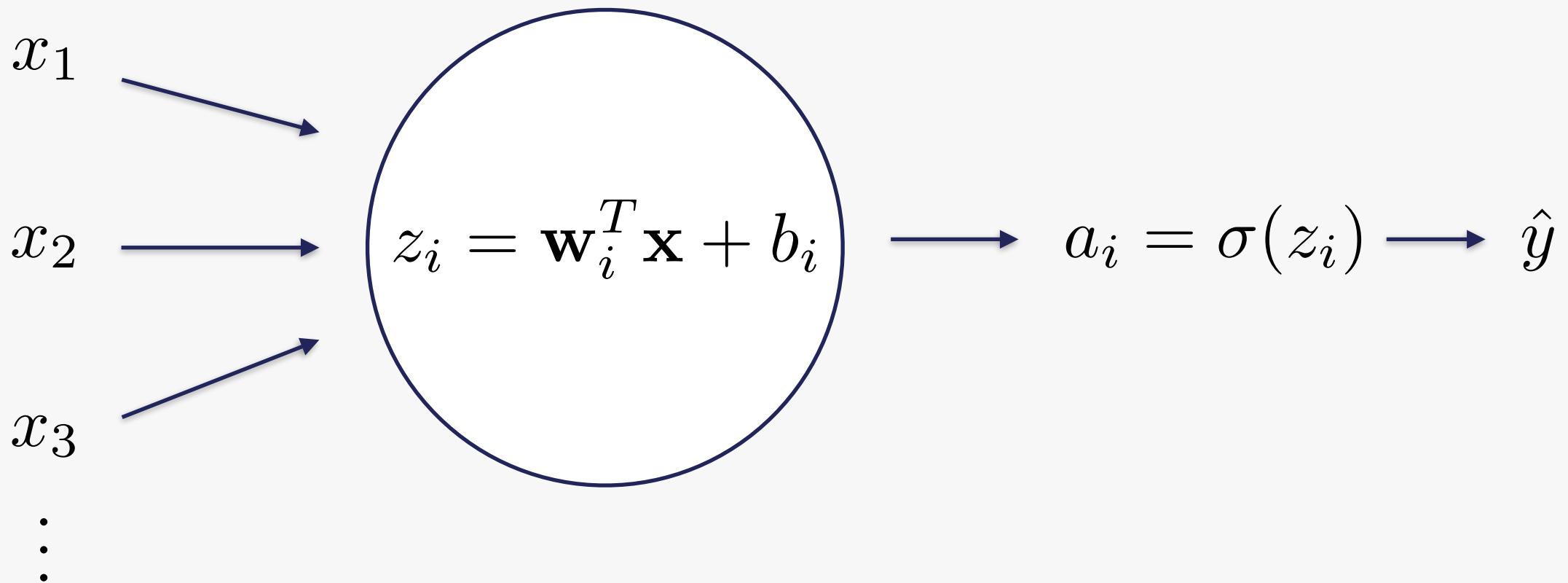
# OLS solved with gradient descent



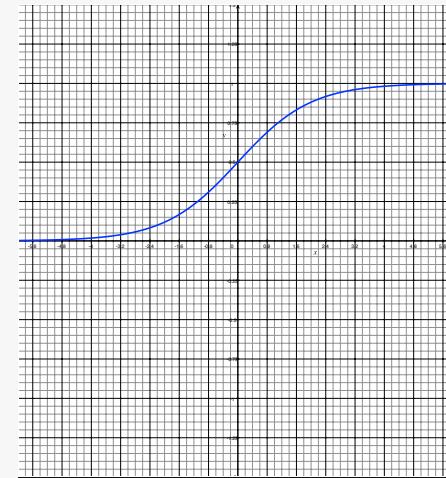
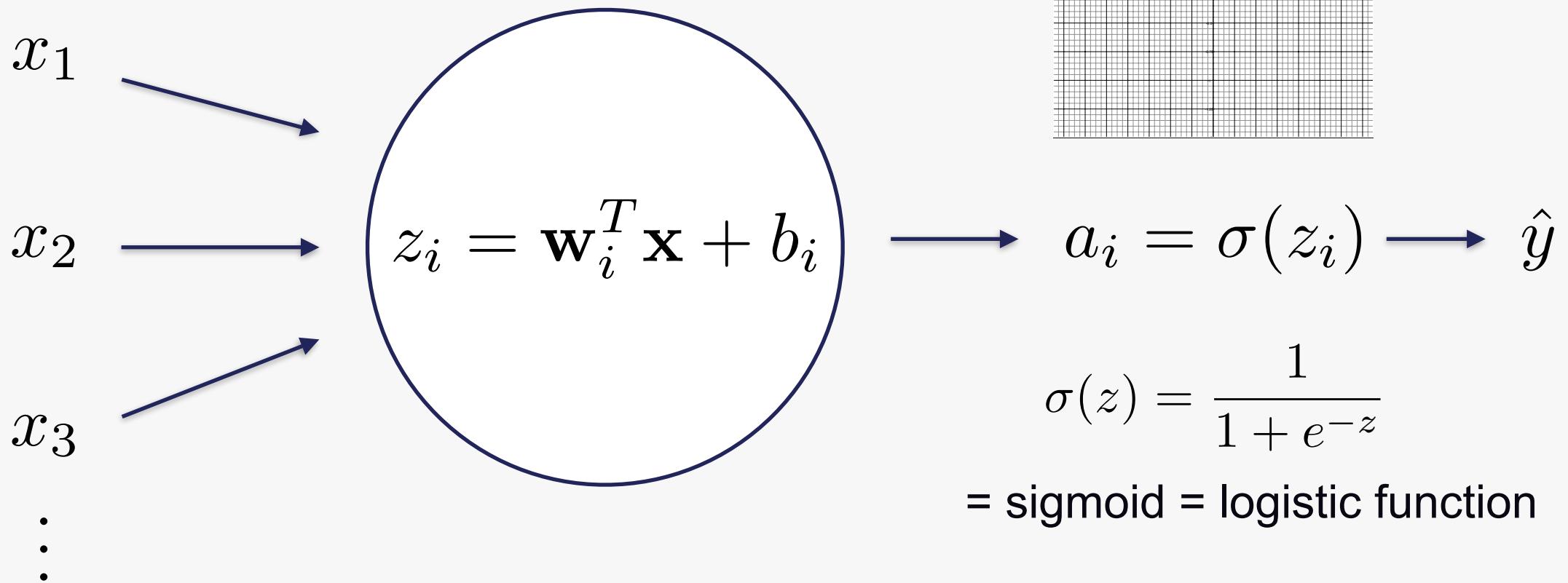
$t = \text{epoch}$   
= neuron sees  
**all samples in the**  
training set

<http://www.briandolhansky.com/blog/artificial-neural-networks-linear-regression-part-1>

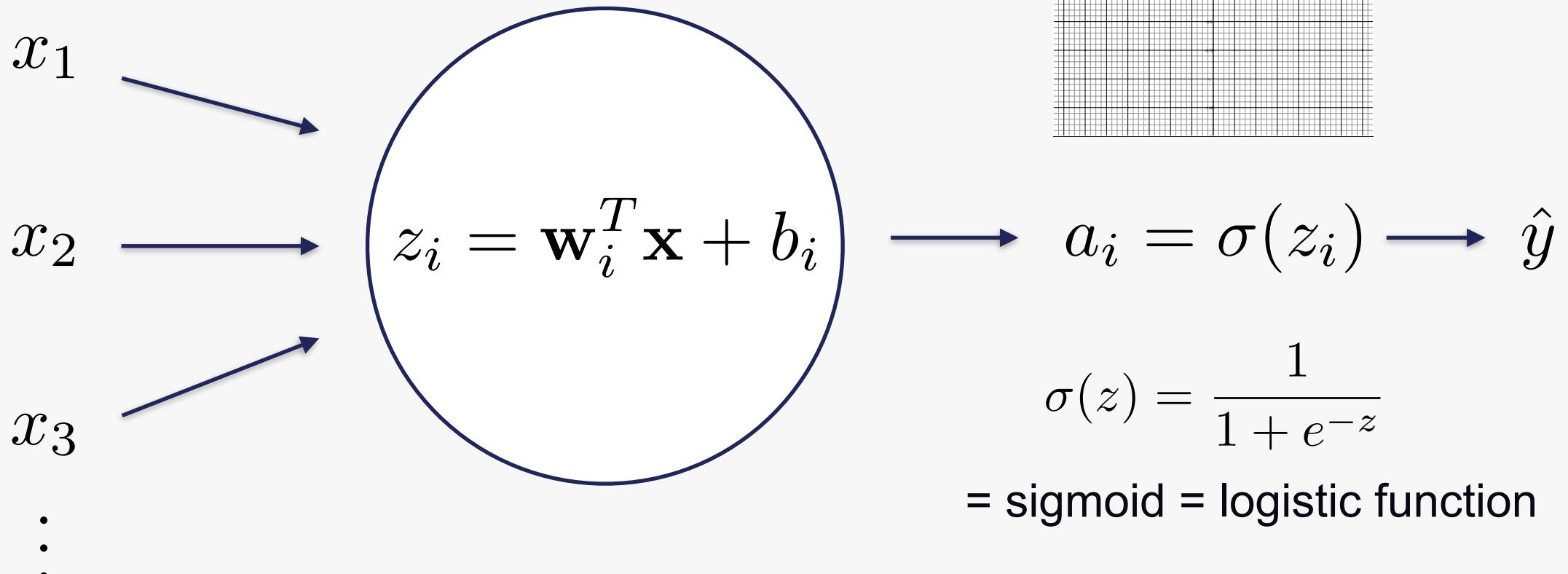
# Neuron (generally)



# Neuron (generally)



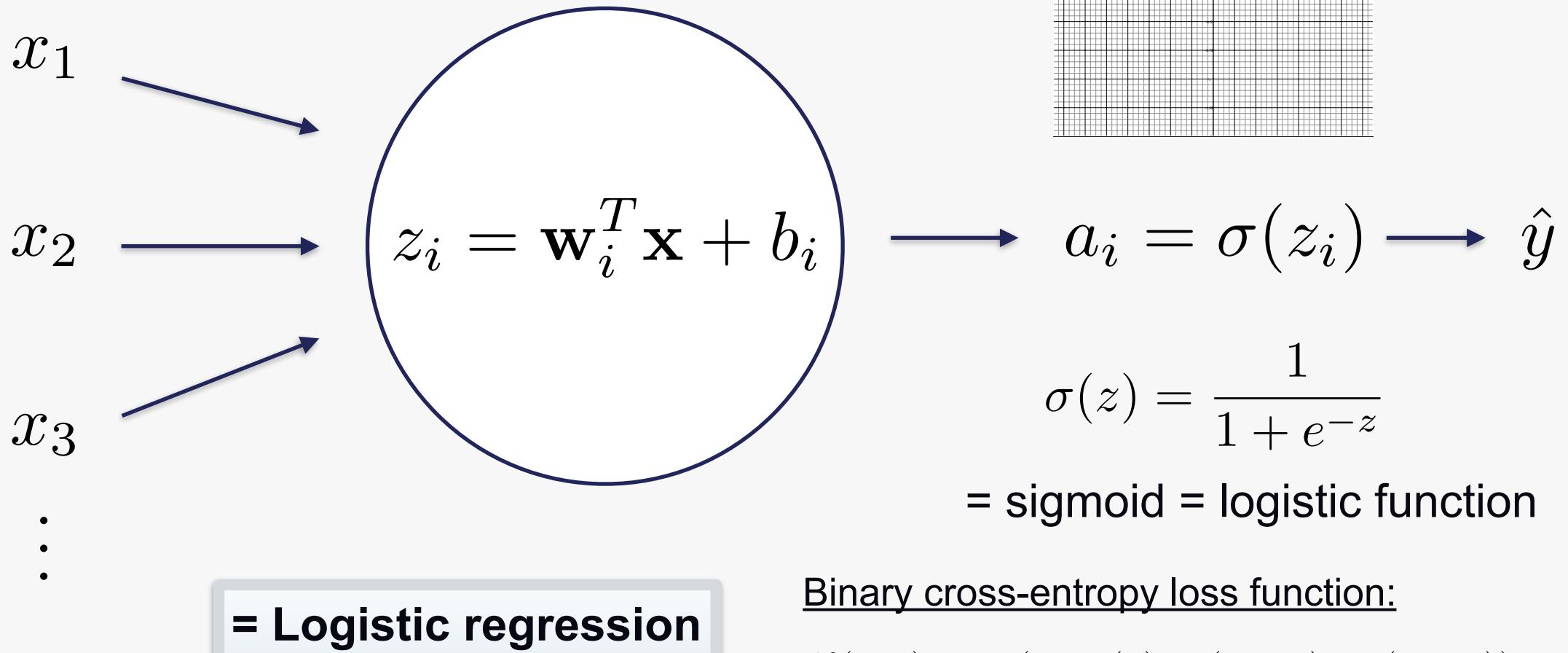
# Neuron (generally)



Binary cross-entropy loss function:

$$\mathcal{L}(\hat{y}, y) = - (y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$$

# Neuron (generally)



# Nonlinear activation functions

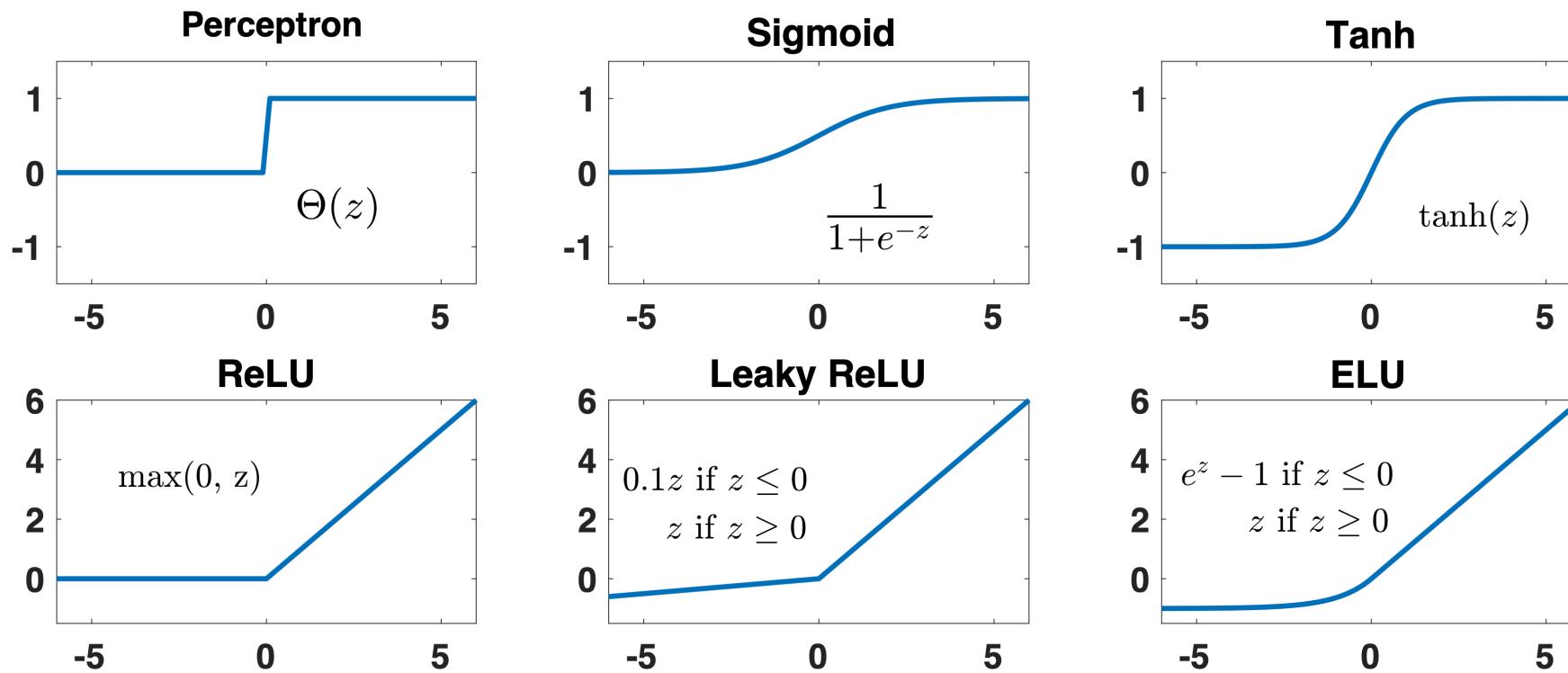
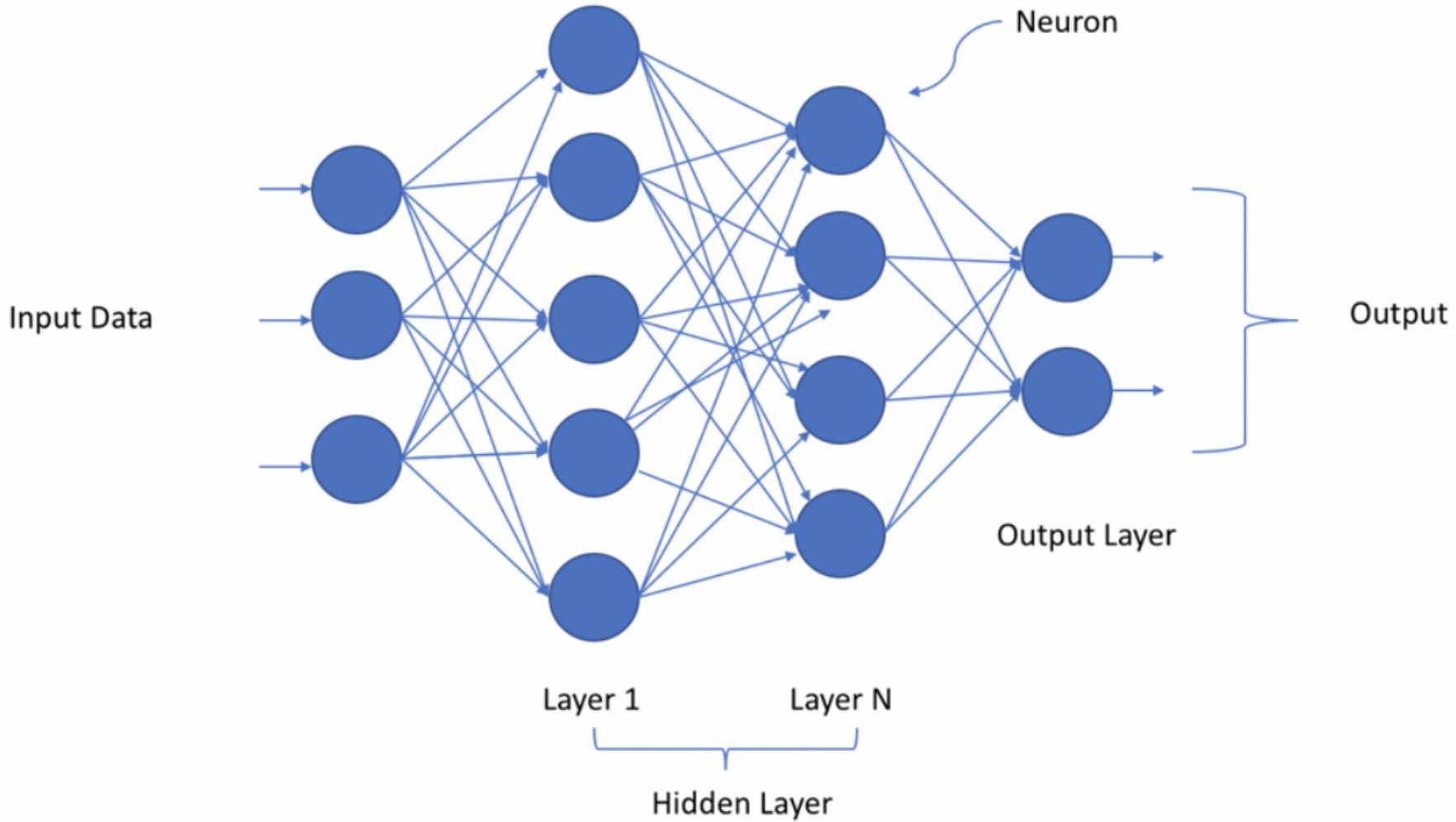


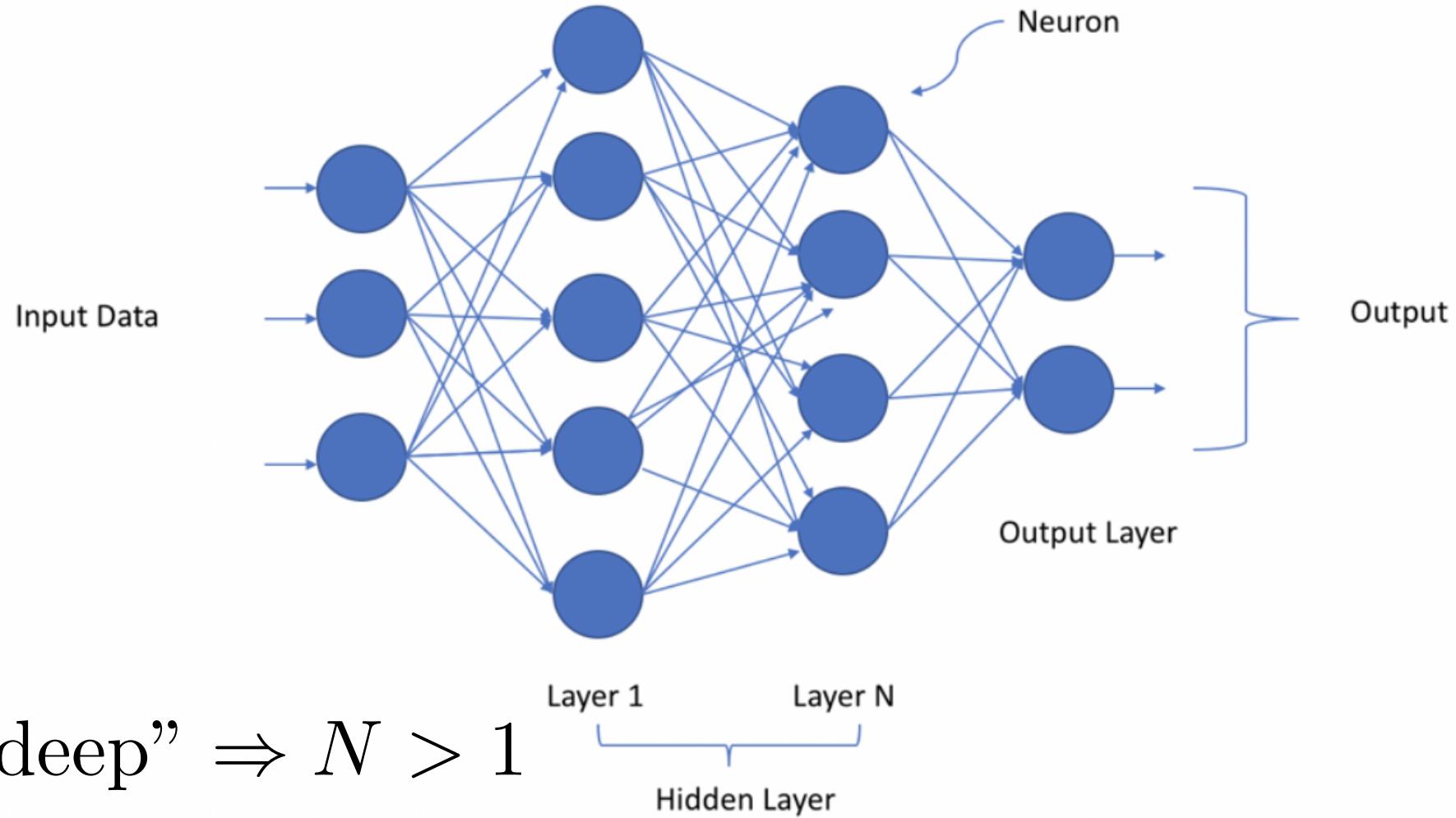
FIG. 36 Possible non-linear activation functions for neurons. In modern DNNs, it has become common to use non-linear functions that do not saturate for large inputs (bottom row) rather than saturating functions (top row).

# Deep neural network (DNN)



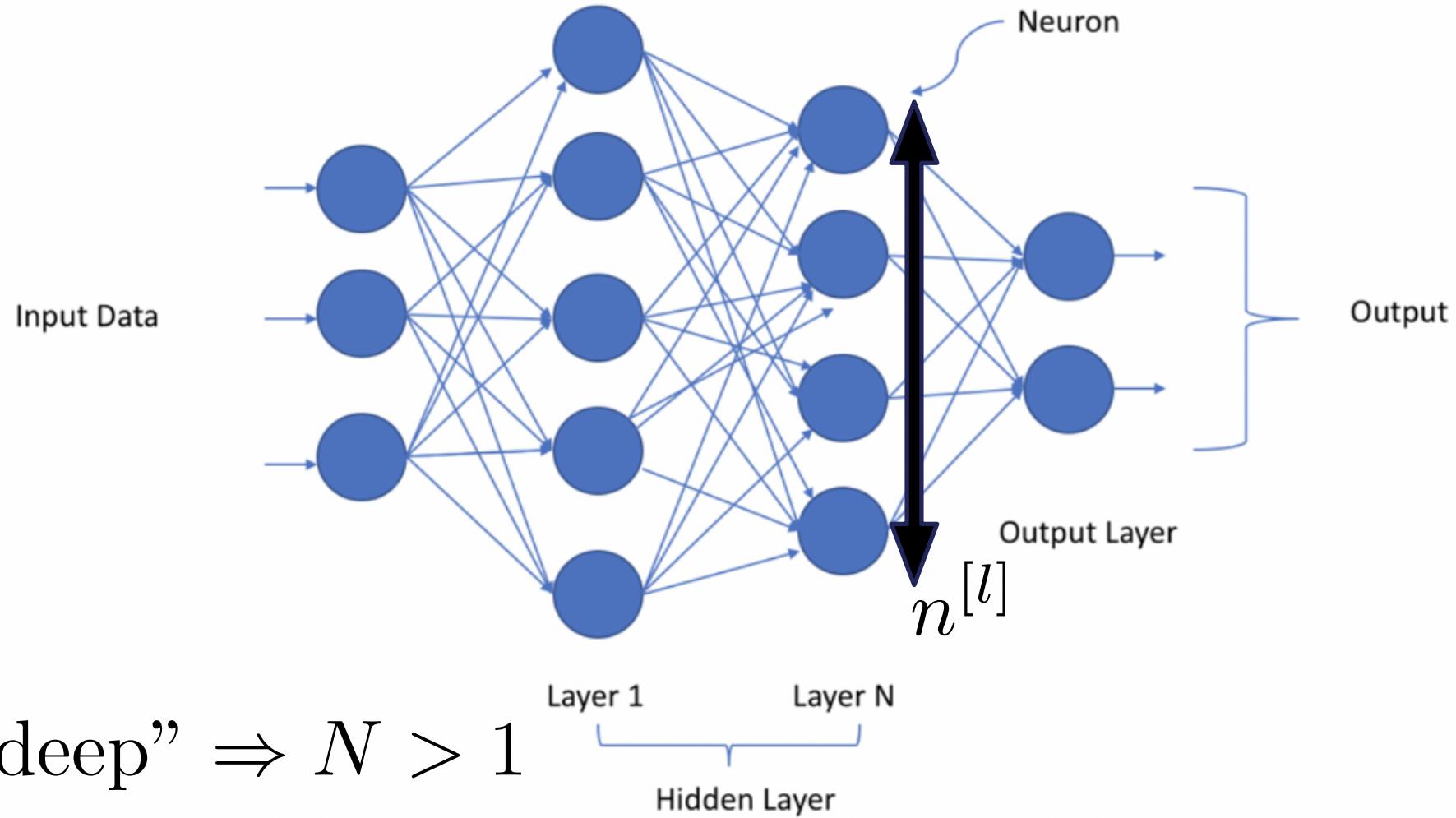
<https://towardsdatascience.com/a-laymans-guide-to-deep-neural-networks-ddcea24847fb>

# Deep neural network (DNN)



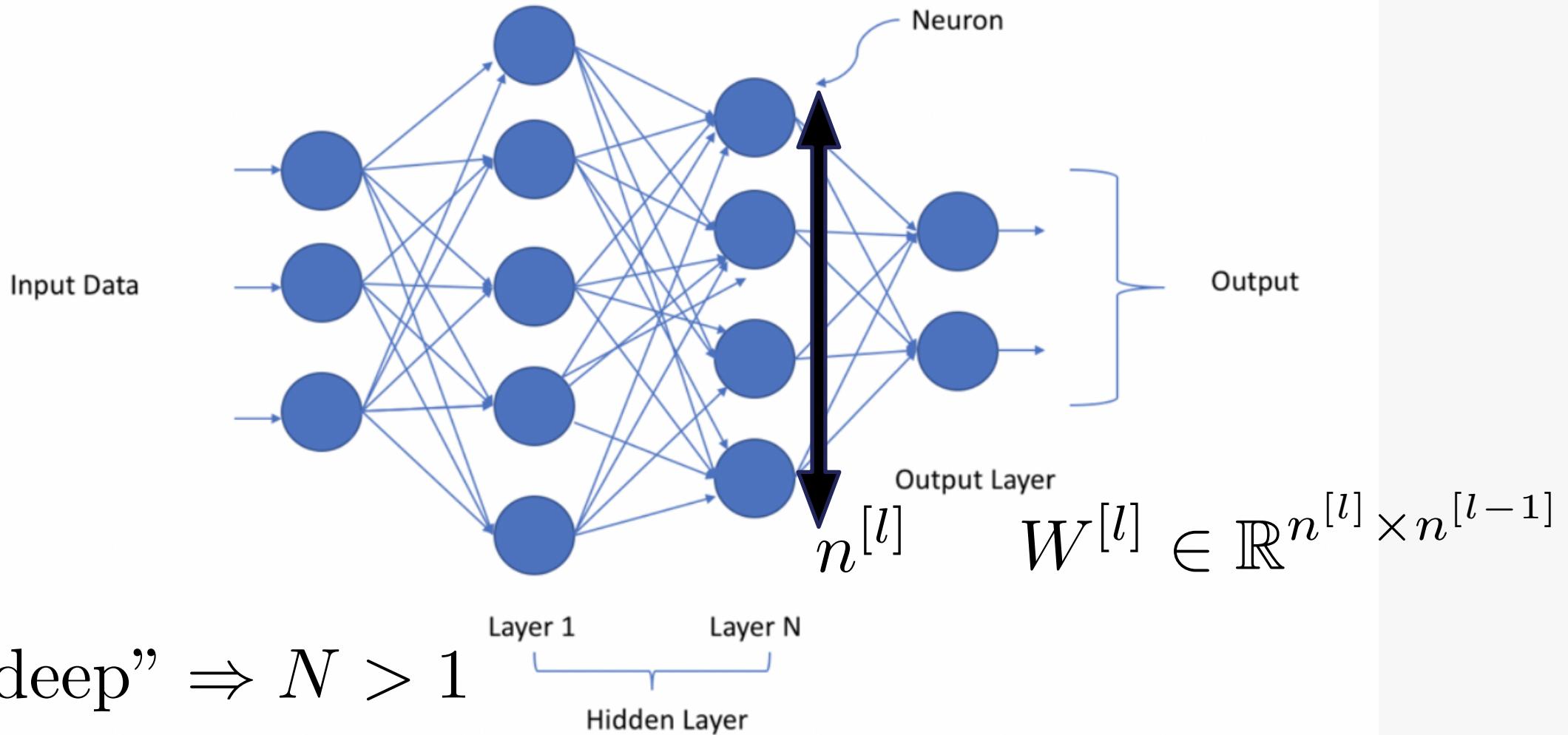
<https://towardsdatascience.com/a-laymans-guide-to-deep-neural-networks-ddcea24847fb>

# Deep neural network (DNN)



<https://towardsdatascience.com/a-laymans-guide-to-deep-neural-networks-ddcea24847fb>

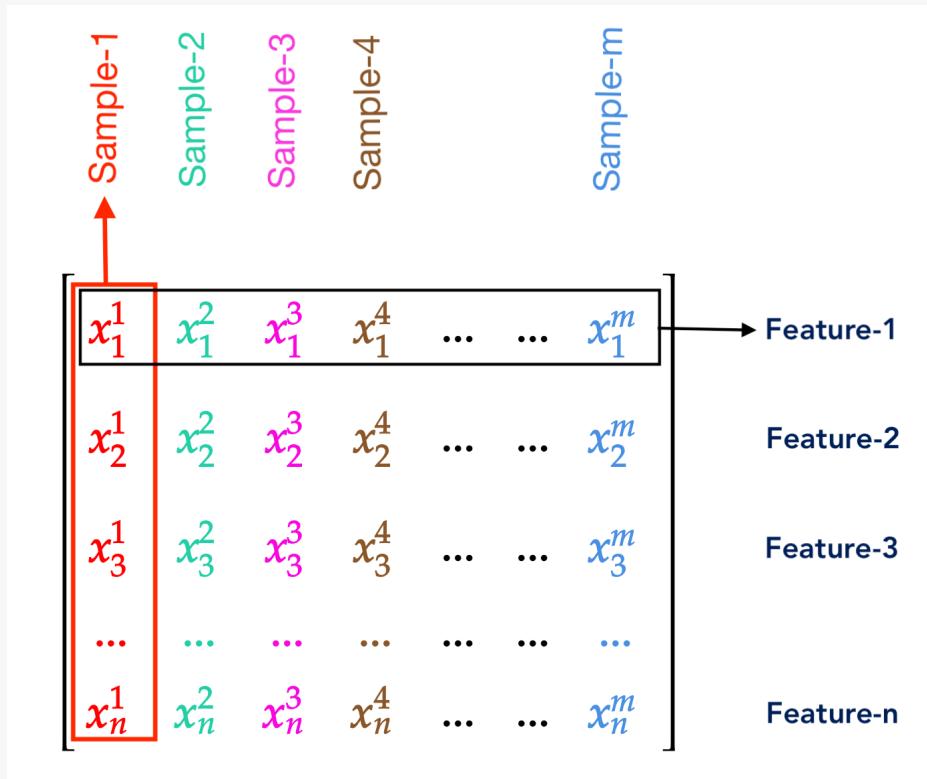
# Deep neural network (DNN)



<https://towardsdatascience.com/a-laymans-guide-to-deep-neural-networks-ddcea24847fb>

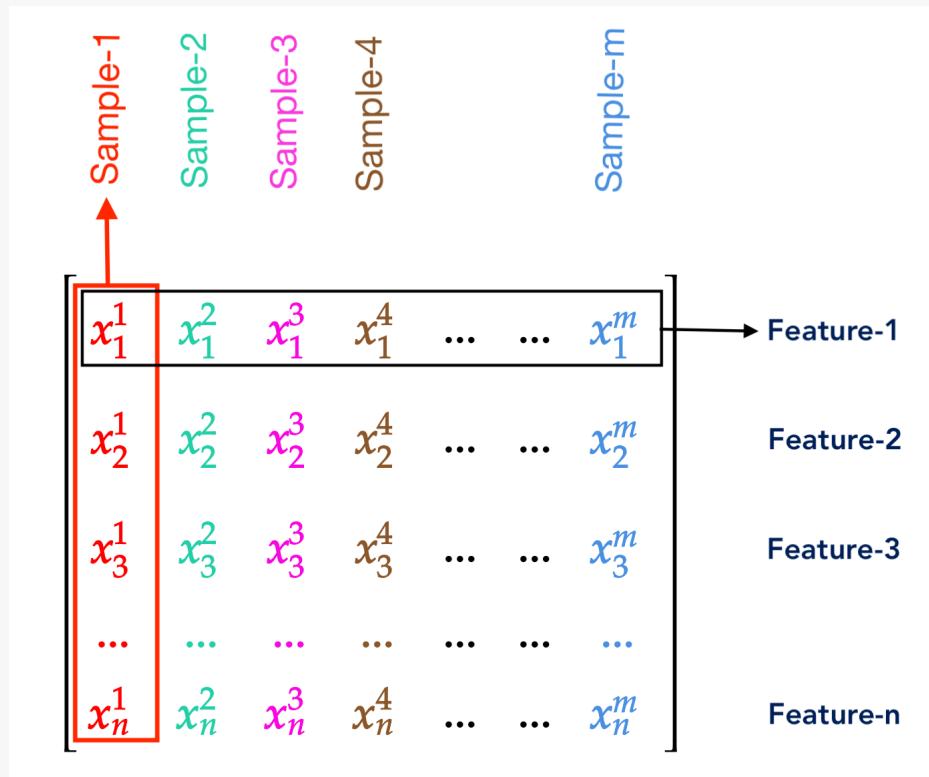
# Vectorized notation: pass $m$ examples at once

Start with horizontally-stacked input  
vectors  $X \in \mathbb{R}^{n \times m}$



# Vectorized notation: pass $m$ examples at once

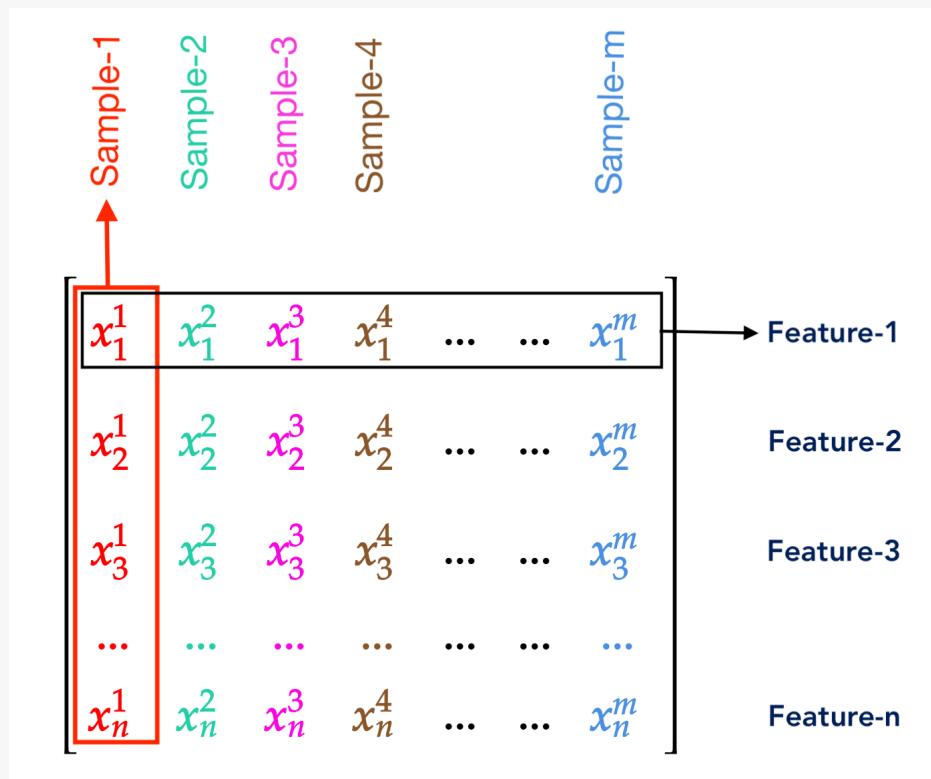
Start with horizontally-stacked input  
vectors  $X \in \mathbb{R}^{n \times m}$



Vectorized forward propagation

# Vectorized notation: pass $m$ examples at once

Start with horizontally-stacked input vectors  $X \in \mathbb{R}^{n \times m}$

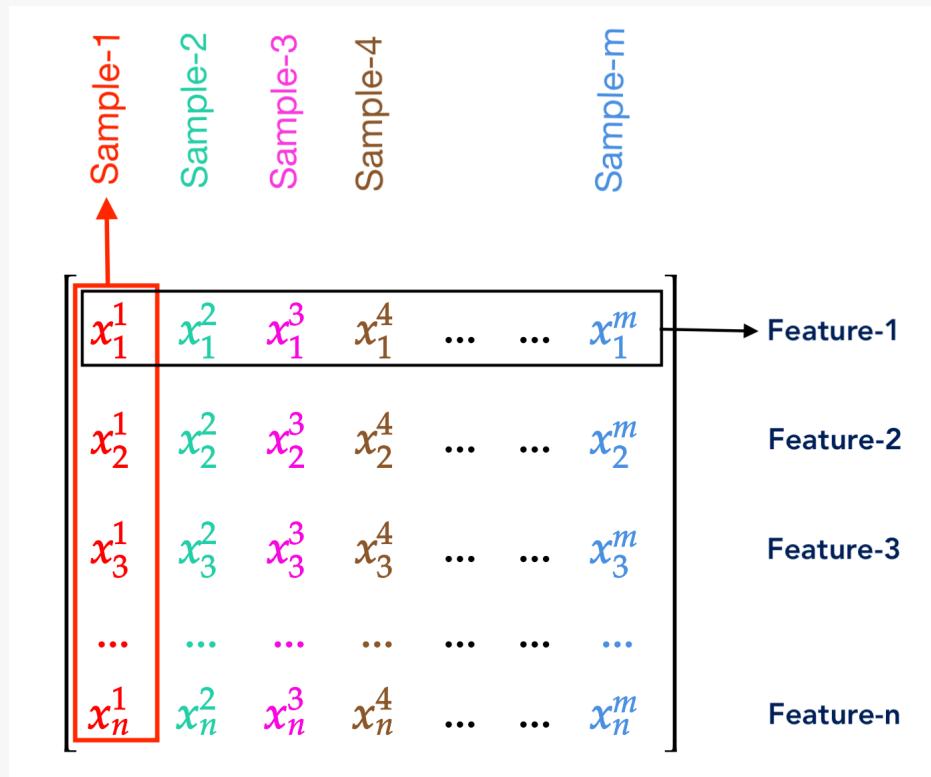


## Vectorized forward propagation

$$Z^{[l]} = W^{[l]} A^{[l-1]} + B^{[l]}$$

# Vectorized notation: pass $m$ examples at once

Start with horizontally-stacked input vectors  $X \in \mathbb{R}^{n \times m}$

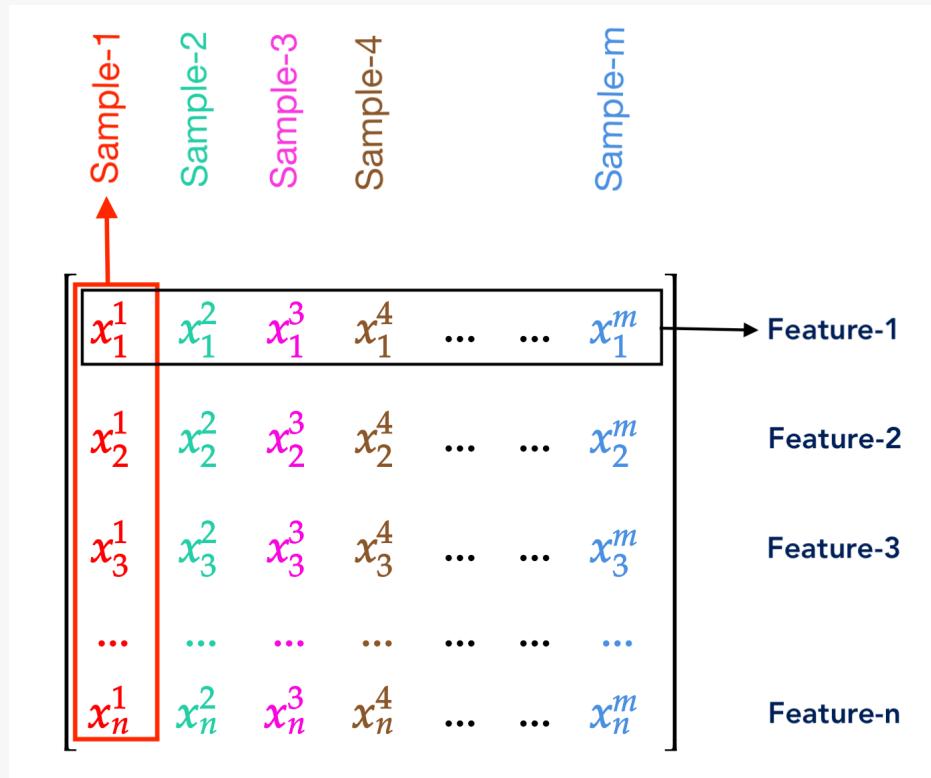


## Vectorized forward propagation

$$Z^{[l]} = W^{[l]} A^{[l-1]} + B^{[l]}$$
$$A^{[l]} = \sigma(Z^{[l]})$$

# Vectorized notation: pass $m$ examples at once

Start with horizontally-stacked input vectors  $X \in \mathbb{R}^{n \times m}$



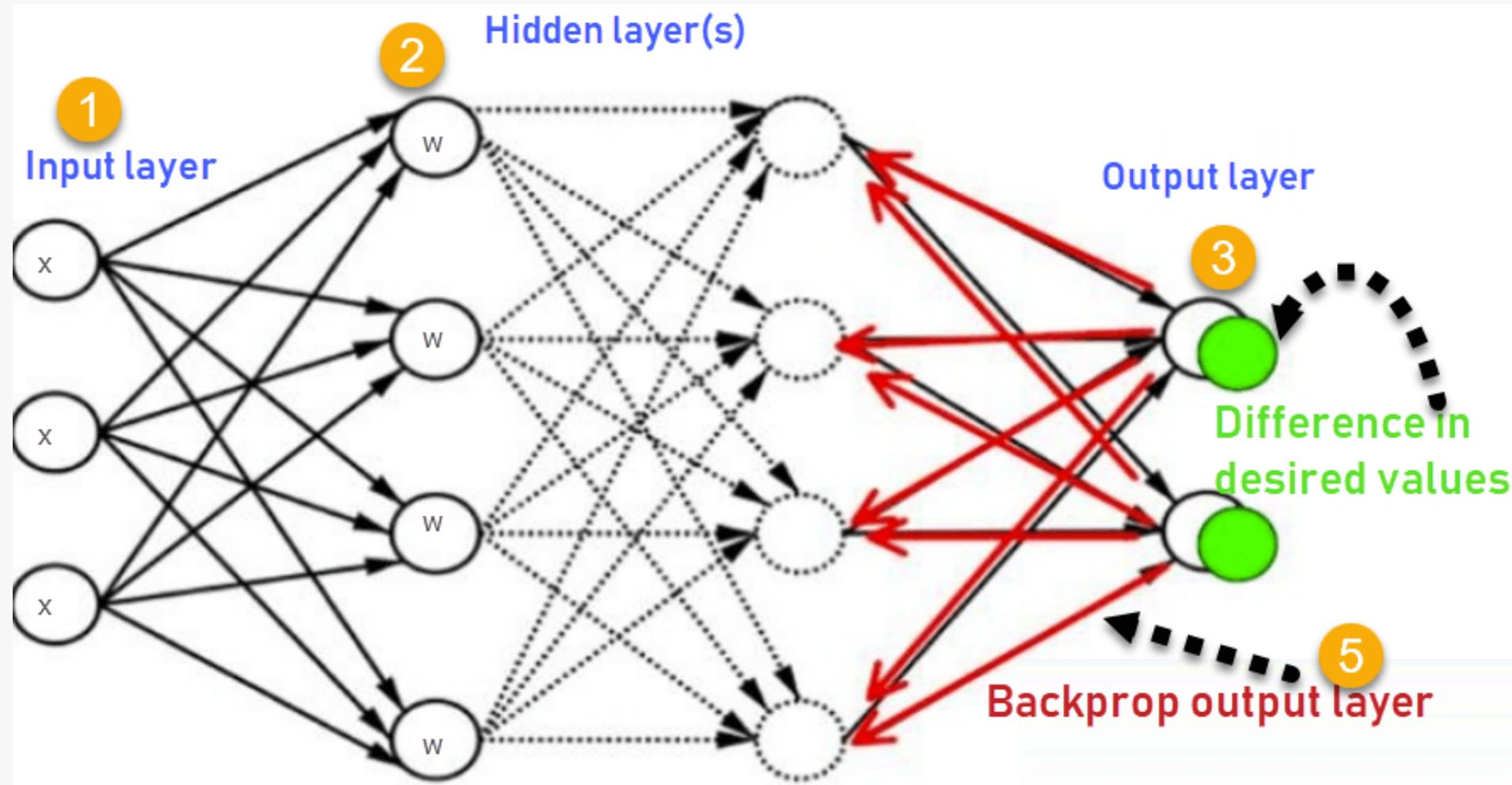
## Vectorized forward propagation

$$Z^{[l]} = W^{[l]} A^{[l-1]} + B^{[l]}$$

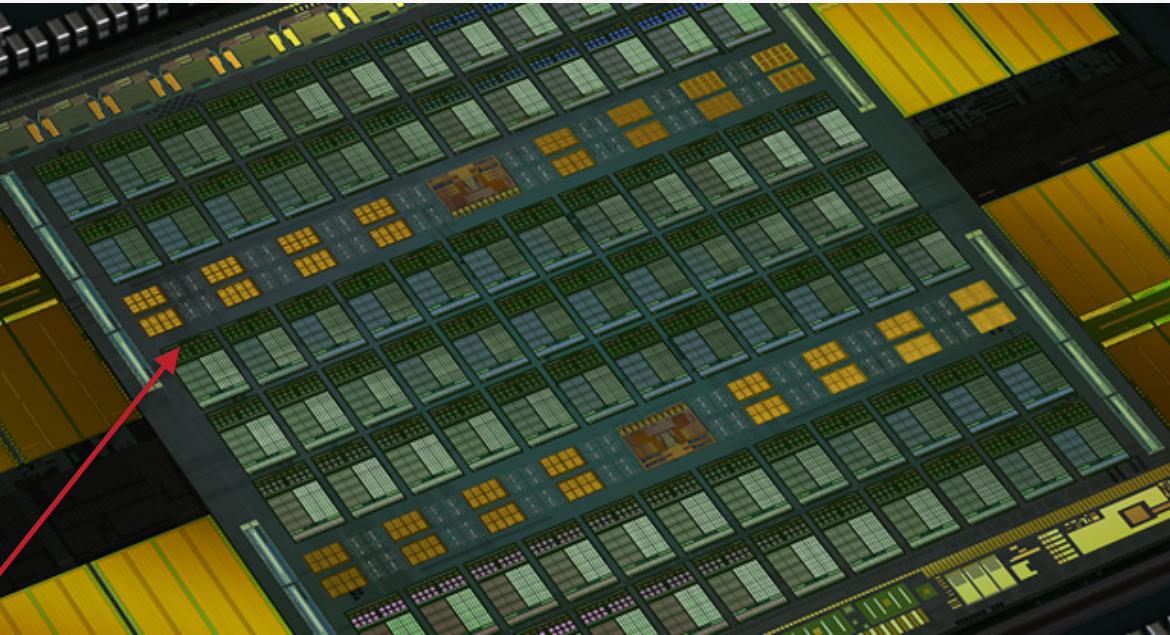
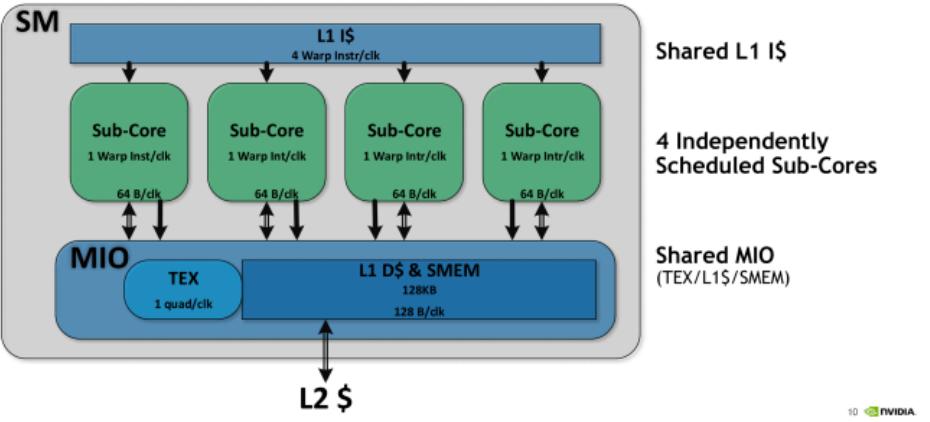
$$A^{[l]} = \sigma(Z^{[l]})$$

$$Z^{[l]}, A^{[l]} \in \mathbb{R}^{n^{[l]} \times m}$$

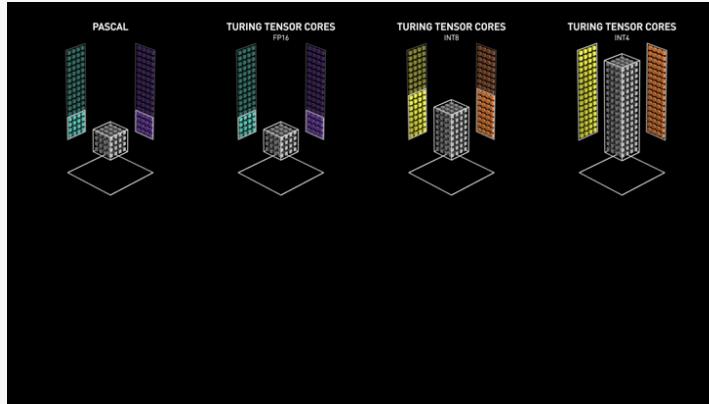
# “Layer-by-layer” gradient descent through backpropagation



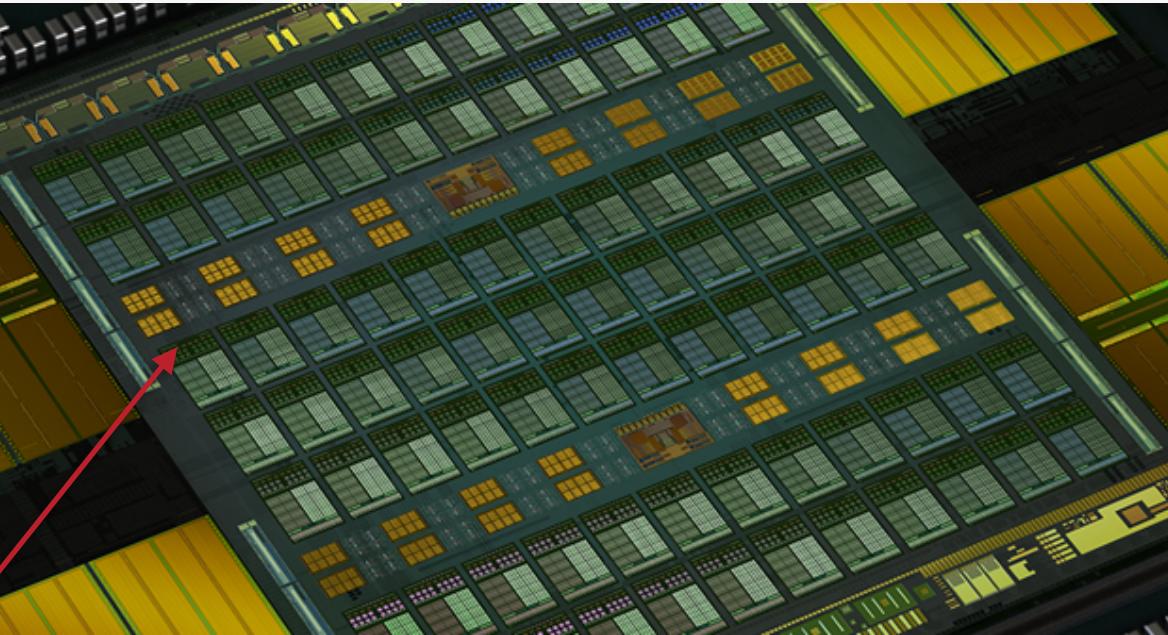
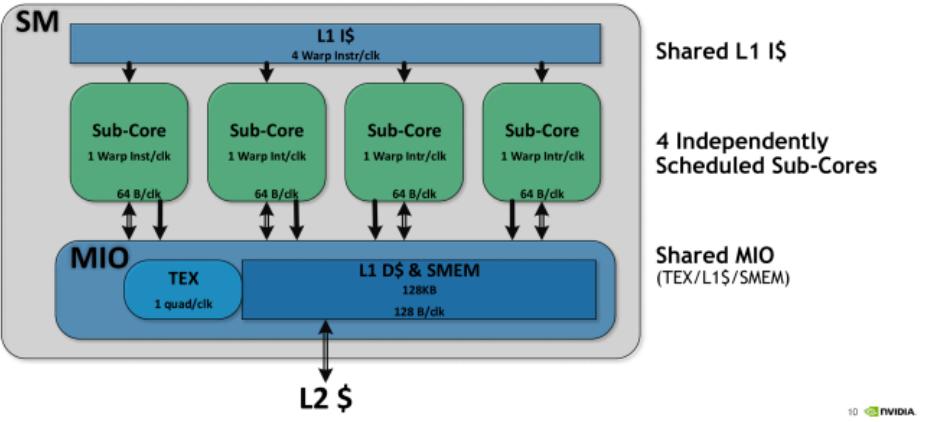
## SM MICROARCHITECTURE



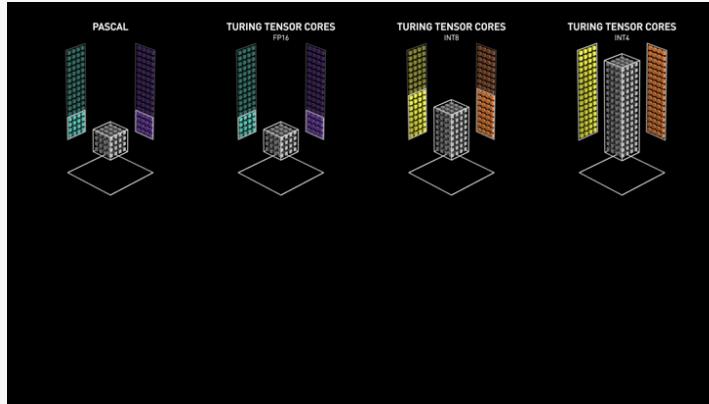
$$\left( \begin{array}{cccc} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{array} \right) + \left( \begin{array}{cccc} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{array} \right) = \left( \begin{array}{cccc} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{array} \right)$$



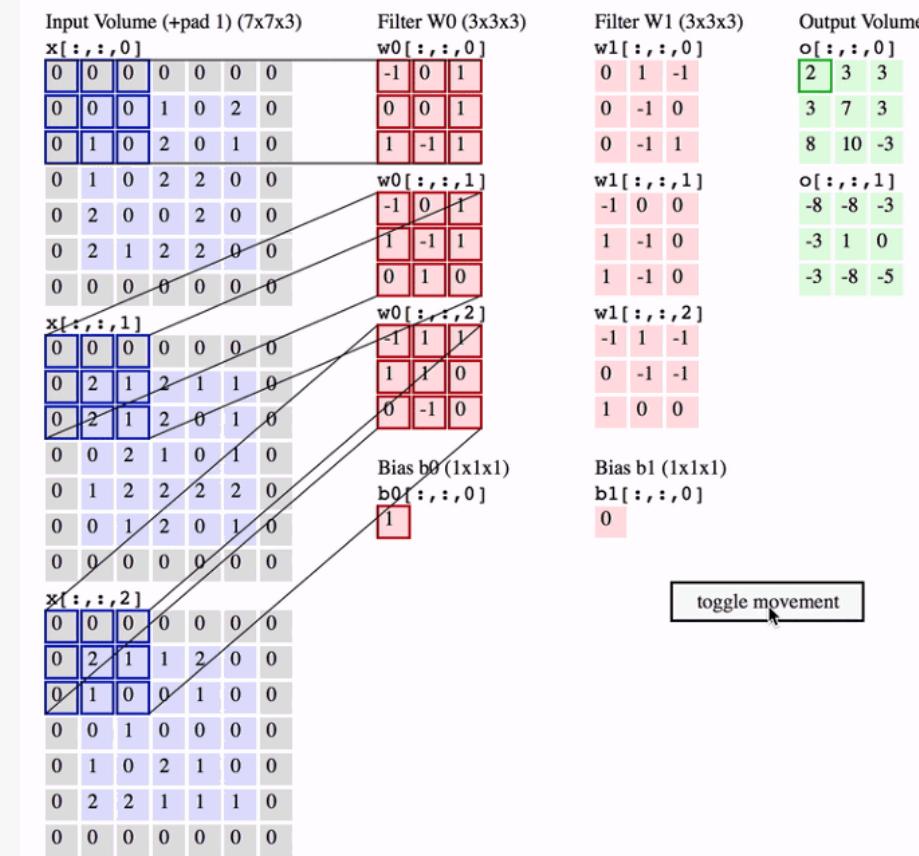
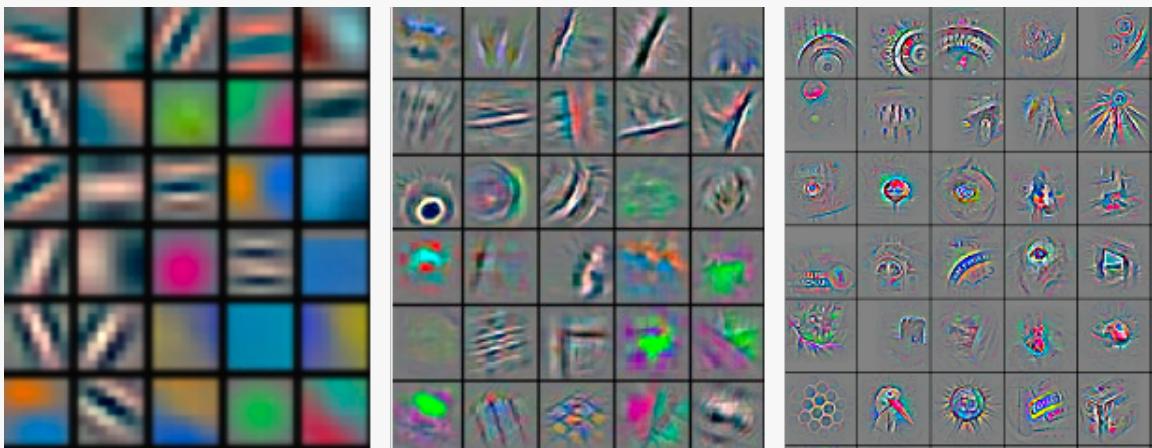
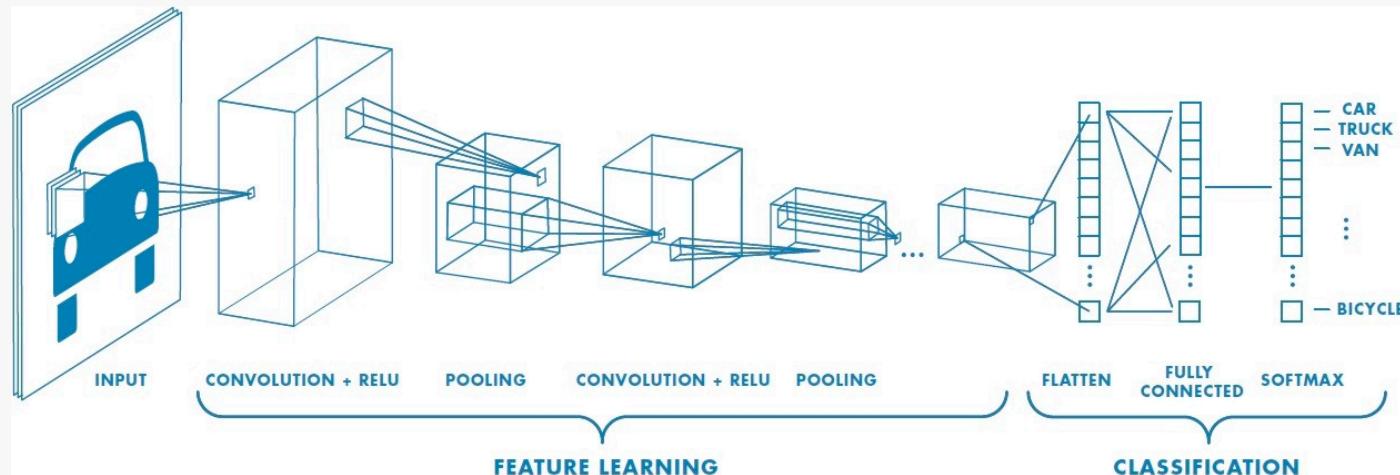
## SM MICROARCHITECTURE



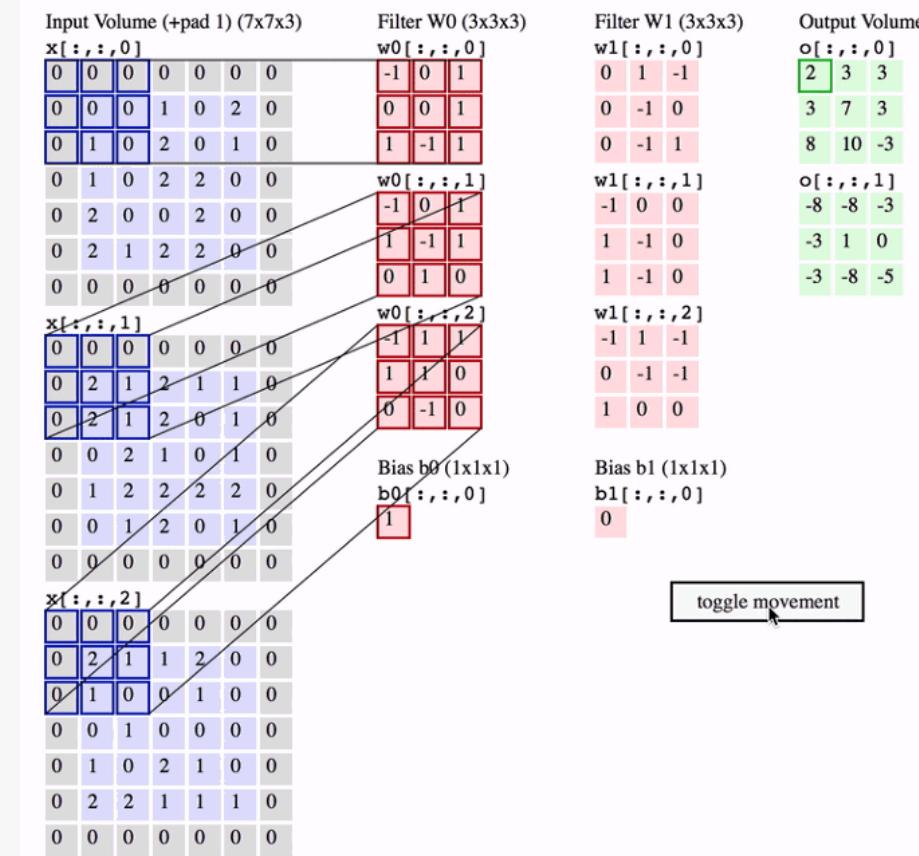
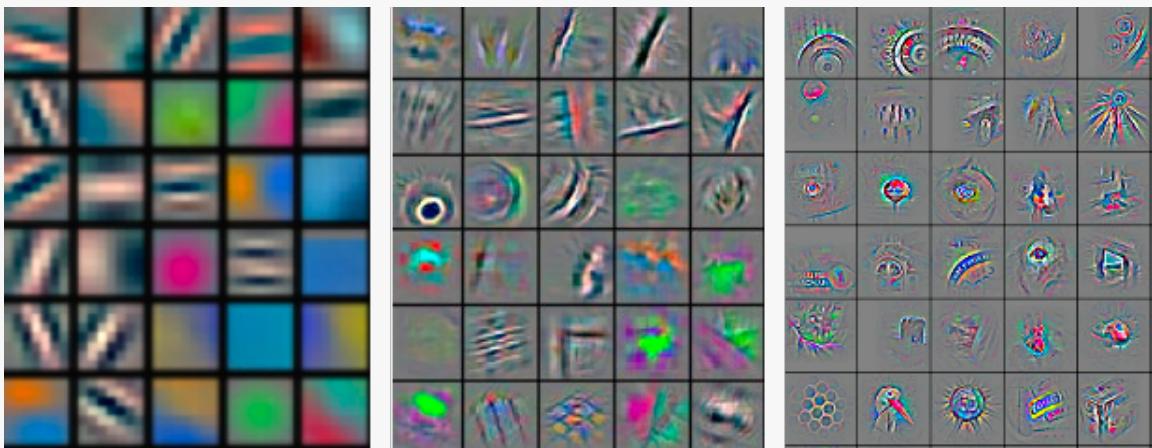
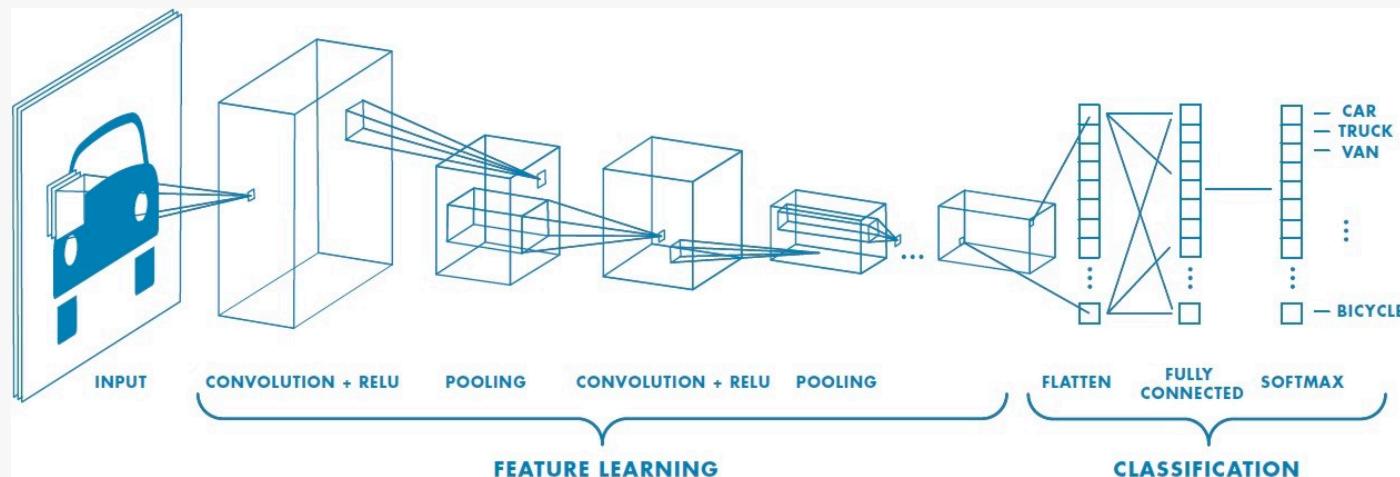
$$\left( \begin{array}{cccc} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{array} \right) + \left( \begin{array}{cccc} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{array} \right) = \left( \begin{array}{cccc} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{array} \right)$$



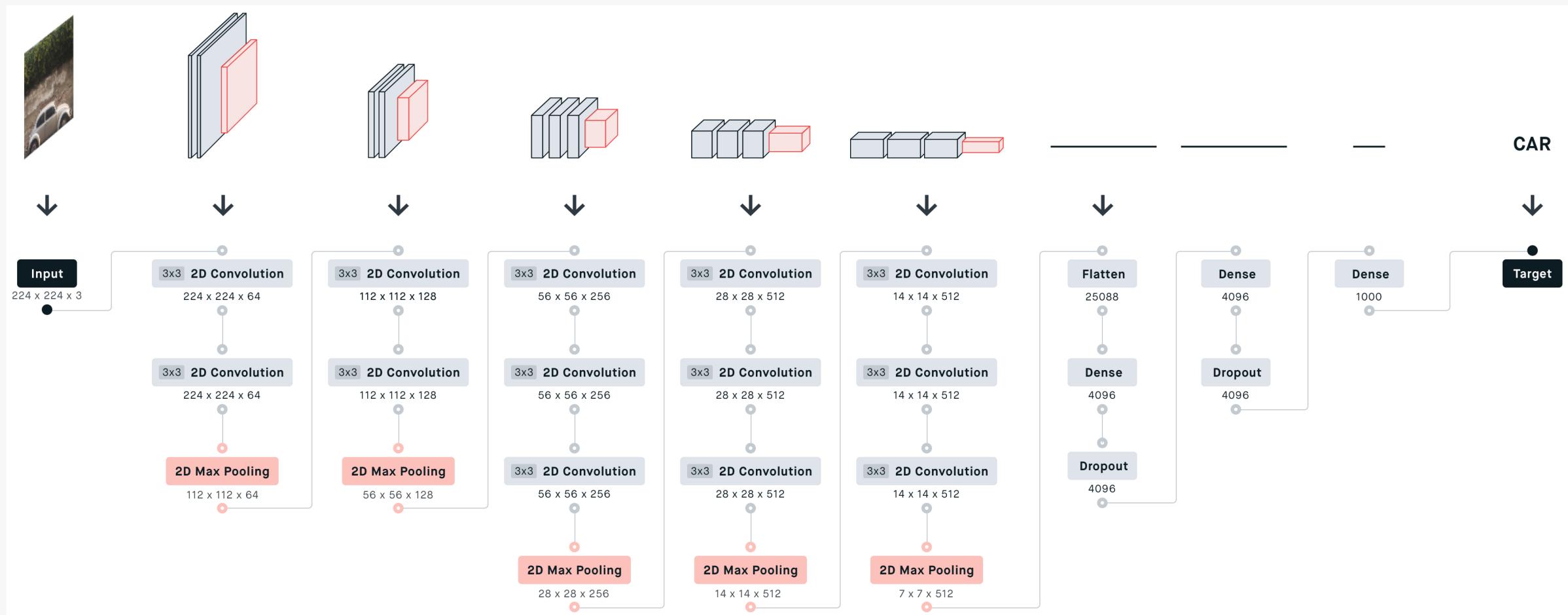
# Convolutional layer



# Convolutional layer

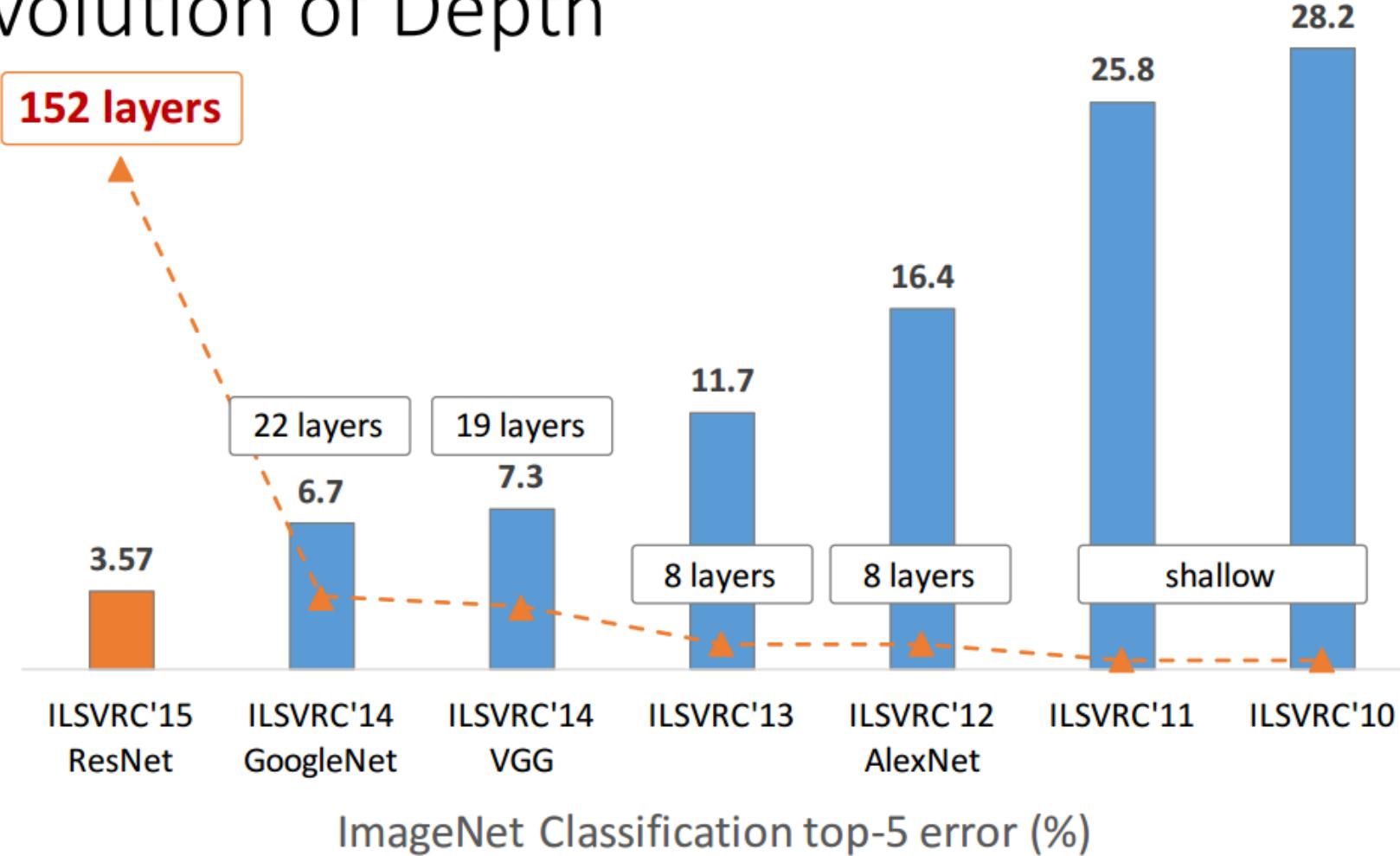


# Convolutional neural network: VGG-16



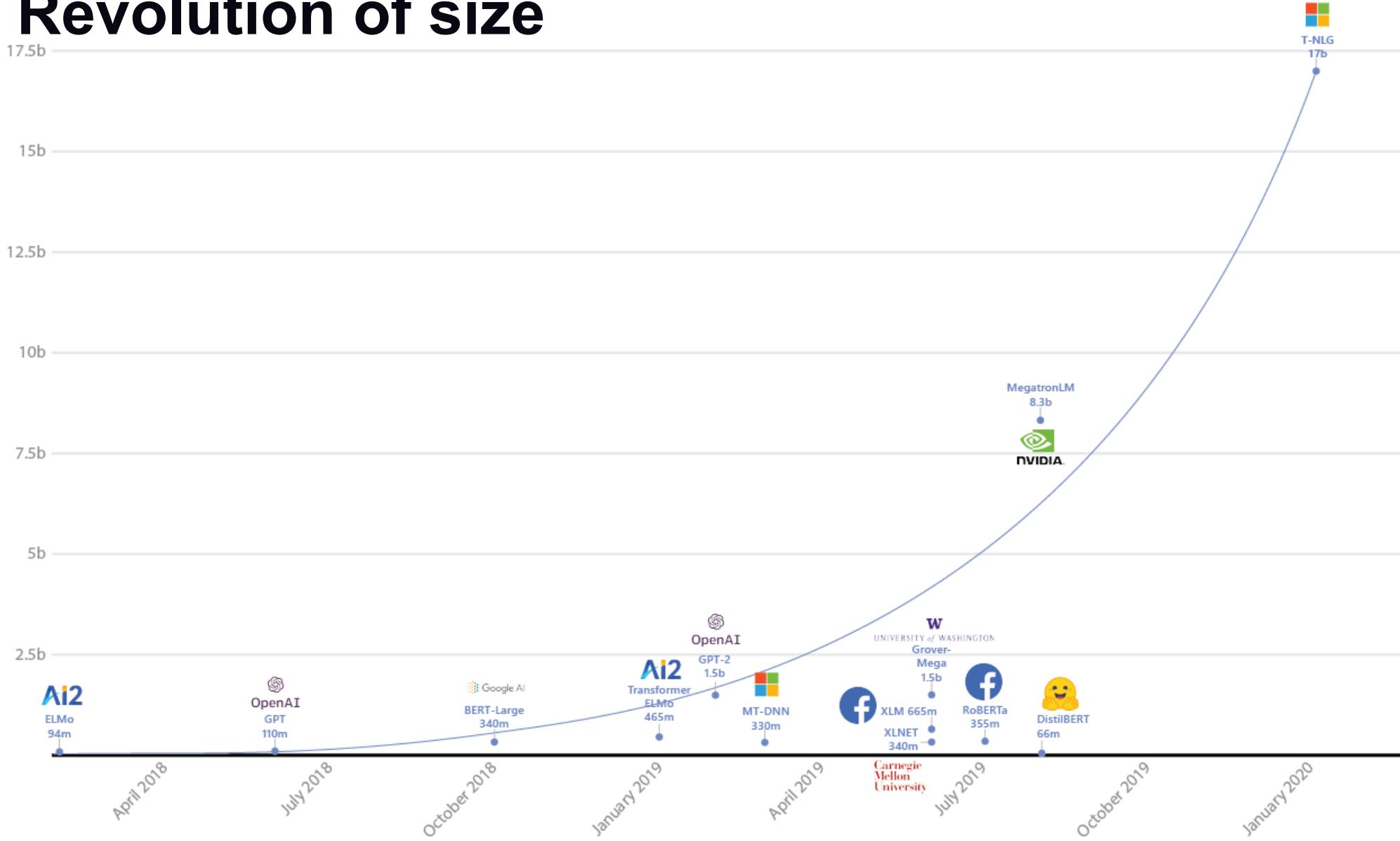
<https://peltarion.com/>

# Revolution of Depth



Kaiming He's ICML 2016 Tutorial

# Revolution of size



<https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/>

# Revolution of size

175b

\* GPT-3  
June 2020

- 
- 
- 

17.5b



17b



# Software frameworks



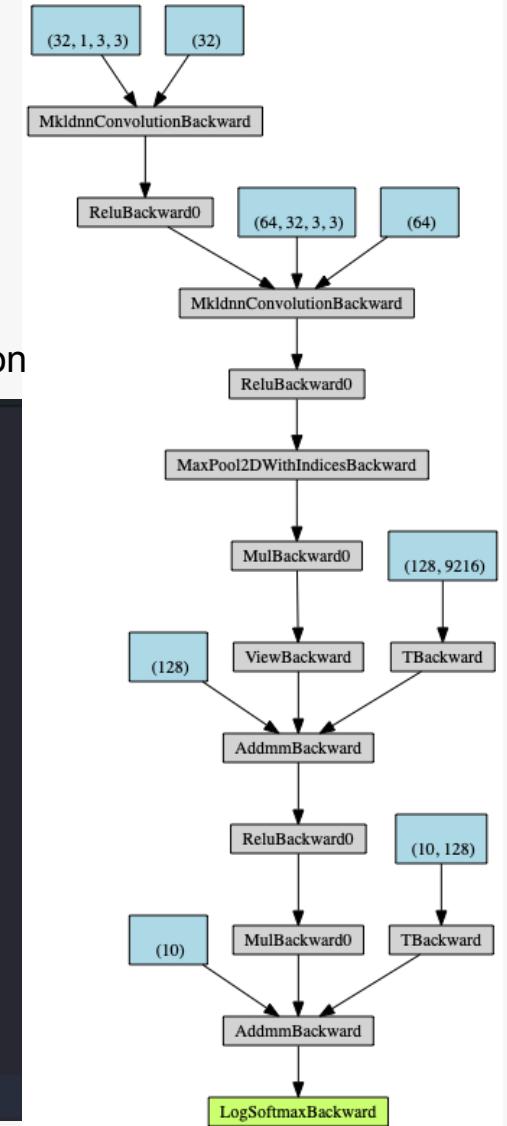
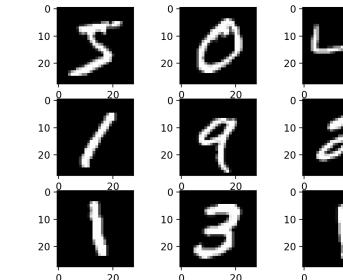
TensorFlow



Keras



PyTorch



Example: PyTorch implementation of 4-layer CNN for MNIST classification

```
class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv1 = nn.Conv2d(1, 10, kernel_size=5)
        self.conv2 = nn.Conv2d(10, 20, kernel_size=5)
        self.conv2_drop = nn.Dropout2d()
        self.fc1 = nn.Linear(320, 50)
        self.fc2 = nn.Linear(50, 10)

    def forward(self, x):
        x = F.relu(F.max_pool2d(self.conv1(x), 2))
        x = F.relu(F.max_pool2d(self.conv2_drop(self.conv2(x)), 2))
        x = x.view(-1, 320)
        x = F.relu(self.fc1(x))
        x = F.dropout(x, training=self.training)
        x = self.fc2(x)
        return F.log_softmax(x)
```

<https://github.com/pytorch/examples/blob/master/mnist/main.py>

# Hands-on resources

## Google Colaboratory

<https://colab.research.google.com/>

## Tutorial Notebooks

[https://github.com/argonne-lcf/ATPESC\\_MachineLearning/](https://github.com/argonne-lcf/ATPESC_MachineLearning/)



# Thank you!